

#3

Docket No. 826.1718

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Patent Application of:)	
)	
Yoshitake SHINKAI, et al.)	
)	Group Art Unit: Unassigned
Serial No.: To be assigned)	
)	Examiner: Unassigned
Filed: March 26, 2001)	
)	
For: FILE REPLICATION SYSTEM,)	
REPLICATION CONTROL)	
METHOD, AND STORAGE)	
MEDIUM)	

31002 U.S. PTO
09/817288
03/27/01

**SUBMISSION OF CERTIFIED COPY OF PRIOR FOREIGN
APPLICATION IN ACCORDANCE
WITH THE REQUIREMENTS OF 37 C.F.R. §1.55**

*Assistant Commissioner for Patents
Washington, D.C. 20231*

Sir:

In accordance with the provisions of 37 C.F.R. §1.55, the applicant submits herewith a certified copy of the following foreign application:

Japanese Patent Application No. 2000-126797
Filed: April 27, 2000.

It is respectfully requested that the applicant be given the benefit of the foreign filing date as evidenced by the certified papers attached hereto, in accordance with the requirements of 35 U.S.C. §119.

Respectfully submitted,

STAAS & HALSEY LLP

Date: March 26, 2001

By: _____

James D. Halsey, Jr.
Registration No. 22,729

700 Eleventh Street, N.W.
Suite 500
Washington, D.C. 20001
(202) 434-1500

PATENT OFFICE
JAPANESE GOVERNMENT

31002 U.S. PRO
09/817288
03/27/01

This is to certify that the annexed is a true copy of the
following application as filed with this Office.

Date of Application: April 27, 2000

Application Number: Patent Application No. 2000-126797

Applicant(s): FUJITSU LIMITED

December 22, 2000

Commissioner,
Patent Office Kozo OIKAWA

Certificate No. 2000-3105863

日 本 国 特 許 庁
PATENT OFFICE
JAPANESE GOVERNMENT



別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日
Date of Application: 2000年 4月27日

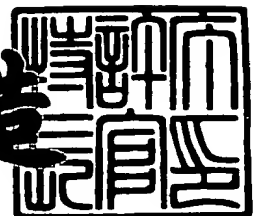
出 願 番 号
Application Number: 特願2000-126797

出 願 人
Applicant(s): 富士通株式会社

2000年12月22日

特 許 庁 長 官
Commissioner,
Patent Office

及 川 耕 造



出証番号 出証特2000-3105863

【書類名】 特許願

【整理番号】 9952124

【提出日】 平成12年 4月27日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 12/00

【発明の名称】 ファイルレプリケーションシステム、ファイルレプリケーション制御方法及び記憶媒体

【請求項の数】 22

【発明者】

 【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

 【氏名】 新開 慶武

【発明者】

 【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

 【氏名】 吉沢 直美

【発明者】

 【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

 【氏名】 塩沢 賢輔

【特許出願人】

 【識別番号】 000005223

 【氏名又は名称】 富士通株式会社

【代理人】

 【識別番号】 100074099

 【住所又は居所】 東京都千代田区二番町8番地20 二番町ビル3F

 【弁理士】

 【氏名又は名称】 大菅 義之

 【電話番号】 03-3238-0031

【選任した代理人】

【識別番号】 100067987

【住所又は居所】 神奈川県横浜市鶴見区北寺尾 7 - 2 5 - 2 8 - 5 0 3

【弁理士】

【氏名又は名称】 久木元 彰

【電話番号】 045-573-3683

【手数料の表示】

【予納台帳番号】 012542

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9705047

【ブルーフの要否】 要

【書類名】 明細書

【発明の名称】 ファイルレプリケーションシステム、ファイルレプリケーション制御方法及び記憶媒体

【特許請求の範囲】

【請求項 1】 複数のノードがネットワークに接続され、該各ノード上に共用ファイルを配置するファイルシステムにおいて、

前記複数のノード内の 1 つである第 1 のノードは、

前記共用ファイルに対するアクセス要求が生じた時、前記複数のノード内の 1 つである第 2 のノードに該共用ファイルに対するアクセスの許可を求める第 1 のトークン管理手段と、

自ノード内で生じた共用ファイルに対するアクセスを受け付け、該アクセス要求に対し前記第 1 のトークン管理手段に前記アクセスの許可獲得を依頼し、該許可が得られない時、前記共用ファイルに対する更新許可を持つノードに該共用ファイルへのアクセス処理を依頼する I O 要求インタセプト手段と、

を備え、

前記第 2 のノードは、

他ノードからの共用ファイルに対するアクセスの許可要求に対し、別のノードに該共用ファイルに対する更新許可を与えている時、該アクセス許可要求に対する応答として該更新許可を与えているノードを通知する第 2 のトークン管理手段を備えることを特徴とするファイルシステム。

【請求項 2】 他のノードとネットワークによって接続され、該他のノードとの共用ファイルを保持するノードにおいて、

前記共用ファイルに対するアクセス要求を管理するトークン管理手段と、

自ノード内で生じた共用ファイルに対するアクセス要求に対し、前記トークン管理手段に該共有ファイルへのアクセス許可を求める I O 要求インタセプト手段と、

を備え、

前記トークン管理手段は、前記 I O 要求インタセプトからのアクセス要求に対し、既に他のノードが前記共用ファイルに対する更新許可を保持する時、該更新

許可を保持するノードを前記 I O 要求インタセプト手段に通知し、前記 I O 要求インタセプト手段は、前記アクセス許可が得られない時、該更新許可を保持するノードに前記共用ファイルへのアクセス処理を依頼することを特徴とするノード。

【請求項 3】 新規参入時に自ノードの保持する共用ファイルのデータの復元処理を行う系構成管理手段を更に備え、前記ファイルの復元処理中に、自ノード内で前記共用ファイルに対するアクセス要求が生じた時、前記 I O 要求インタセプト手段は、前記共用ファイルを共用している他のノードにアクセス処理を依頼することを特徴とする請求項 2 に記載のノード。

【請求項 4】 前記共用ファイルへの更新時に更新内容を他の更新との依存関係を示す情報と共に他のノードへ伝播する変更データ通知手段と、

前記依存関係を示す情報に基づいて、更新の順序性を保証しつつ前記更新内容を前記共用ファイルに反映させる受信データ処理手段を更に備えることを特徴とする請求項 2 又は 3 に記載のノード。

【請求項 5】 1 乃至複数の共用ファイル毎に更新内容の伝播方式についての情報を保持する系状態情報保持手段を更に備え、前記変更データ通知手段は、前記形状情報保持手段内の情報に基づいて前記更新内容を伝播することを特徴とする請求項 4 に記載のノード。

【請求項 6】 前記伝播方式は前記共用ファイルを共用する全てのノードに前記更新内容が伝播されるのを保証する同期方式、前記前記共用ファイルを共用する半数のノードに前記更新内容が伝播されるのを保証する半同期方式、及び前記共用ファイルを共用するノードへの前記更新内容の伝播を確認しない非同期方式のいずれか 1 つであることを特徴とする請求項 5 に記載のノード。

【請求項 7】 前記系状態情報保持手段は、前記 1 乃至複数の共用ファイル毎に該共用ファイルを共用するノードについての情報情報をも保持することを特徴とする請求項 4 乃至 6 のいずれか 1 に記載のノード。

【請求項 8】 複数のノードがネットワークに接続され、該ノードが共用ファイルを共用する構成のシステムにおけるファイルレプリケーション制御方法であって、

前記共用ファイルに対するアクセスを行うアクセス要求ノードは、ノードが前記共用ファイルに対する最新のデータを自己が保持する時、自己の共用ファイルにアクセスし、

前記最新のデータを他ノードが保持する時、前記共用ファイルに対するアクセスを該他ノードに依頼することを特徴とするファイルレプリケーション制御方法

【請求項 9】 前記共用ファイルへの更新許可は 1 つのノードにのみ与えられ、前記アクセス要求ノードは共用ファイルにアクセスする時に、他ノードが前記共用ファイルへの更新許可を保持している時、該更新許可を保持しているノードに前記共用ファイルへのアクセス処理を依頼することを特徴とする請求項 8 に記載のファイルレプリケーション制御方法。

【請求項 10】 前記共用ファイルへの更新を行ったノードは、更新内容を他ノードに非同期で伝播し、

前記更新内容が伝播中に他ノードで生じた共用ファイルへのアクセス要求を前記更新を行ったノードが処理することを特徴とする請求項 8 又は 9 に記載のファイルレプリケーション制御方法。

【請求項 11】 前記更新許可を保持しているノードは、自己の更新が依存する更新が全ノードに伝わった後、該更新許可の解放を行うことを特徴とする請求項乃至 10 のいずれか 1 つに記載のファイルレプリケーション制御方法。

【請求項 12】 前記共用ファイルへの更新内容は順序性を保証して反映されることを特徴とする請求項 8 乃至 11 のいずれか 1 つに記載のファイルレプリケーション制御方法。

【請求項 13】 他の更新との順序関係を示す依存情報を前記更新内容と共に他ノードに伝播することを特徴とする請求項 12 に記載のファイルレプリケーション制御方法。

【請求項 14】 前記更新内容を受信したノードは、前記依存情報に基づき、該更新内容に先行する更新内容を受信した後で、該更新内容を自己の共用ファイルへ反映させることを特徴とする請求項 13 に記載のファイルレプリケーション制御方法。

【請求項 1 5】 前記共用ファイルへの更新内容の他ノードへの伝播の方式を 1 乃至複数の前記共有ファイル単位で指定することを特徴とする請求項 9 乃至 1 4 のいずれか 1 つに記載のファイルレプリケーション制御方法。

【請求項 1 6】 前記共用ファイルへの更新内容を伝播するノードを 1 乃至複数の前記共有ファイル単位で指定することを特徴とする請求項 9 乃至 1 5 のいずれか 1 つに記載のファイルレプリケーション制御方法。

【請求項 1 7】 新規参入時に自ノードの保持する共用ファイルのデータの復元処理を行い、該復元処理完了前にユーザプログラムを稼動させることを特徴とする請求項 9 乃至 1 6 のいずれか 1 つに記載のファイルレプリケーション制御方法。

【請求項 1 8】 前記復元処理によるデータの送信は、前記共用ファイルへの更新要求に対する処理と順序性を保証して行われることを特徴とする請求項 1 7 に記載のファイルレプリケーション制御方法。

【請求項 1 9】 前記復元処理完了前に生じた前記共用ファイルへのアクセス要求に対する処理を、前記共用ファイルを共用している他のノードに依頼することを特徴とする請求項 1 7 又は 1 8 に記載のファイルレプリケーション制御方法。

【請求項 2 0】 共用ファイルに対する処理を、該共用ファイルを共用する他ノードと同期して停止する整然停止を行ったとき、整然停止を行ったノードは該整然停止を行ったことを記憶し、該共用ファイルへの処理を再開する際、他ノードと同期して再開することにより該共用ファイルに対するデータを復元処理を行わないことを特徴とする請求項 9 乃至 1 9 に記載のファイルレプリケーション制御方法。

【請求項 2 1】 複数のノードがネットワークに接続される構成のシステムにおけるファイルレプリケーション方法であって、

第 1 のノードはファイルにアクセスする時に、トークン獲得を要求し、

前記要求に対し前記第 1 のノードがトークンを獲得できない時は該トークンを

保持している第 2 のノードを前記第 1 のノードに通知し、

前記第 1 のノードは、前記獲得できない事を通知された時、前記第 2 のノードに前記ファイルアクセスを依頼する

ことを特徴とするファイルレプリケーション方法。

【請求項 2 2】 ネットワークにより他ノードと接続されるノードを構成するコンピュータにより使用された時、

前記各ノードが共用する共用ファイルに対するアクセスを行う時、前記共用ファイルに対する最新のデータを自己が保持するときは自己の共用ファイルにアクセスし、

前記最新のデータを他ノードが保持するときは、前記共用ファイルに対するアクセスを該他ノードに依頼することを前記コンピュータに行わせるためのプログラムを記憶した前記コンピュータが読み出し可能な記録媒体。

【発明の詳細な説明】

【 0 0 0 1 】

【発明の属する技術分野】

本発明は、複数のコンピュータ間にファイルの複製を動的に配置し、負荷分散をはかり性能向上を実現すると共に信頼性を向上させるファイルレプリケーション技術に関する。

【 0 0 0 2 】

【従来の技術】

従来、ネットワークで接続されている複数の計算機システム（ノード）間に同一のデータを動的に配置し、信頼性を向上させる方式として、ファイルレプリケーション技術が知られている。

【 0 0 0 3 】

ファイルレプリケーションでは、あるノード上のファイルが更新されたタイミングでファイルの更新内容を知り、予め定義された他のノード群に変更されたデータのみを伝播してファイルを更新させる。

【 0 0 0 4 】

更新内容の伝播の仕方としては、更新命令の完了がユーザプログラムに通知さ

れた時点で、他ノードへの伝播が完了していることを保証する同期型と、システム内に変更内容を蓄積し、適当なタイミングで他ノードへ伝播する非同期型の伝播が知られている。このうち非同期型の場合は、応答性が良く処理性能は高いが、更新命令の完了がユーザプログラムに通知された時点で、他ノードに更新内容が伝わっている保証はない。

【 0 0 0 5 】

一方従来のファイルレプリケーション方式では、各ノードが保持するデータの同一性あるいは一貫性が保証されていないため以下の問題が発生する。

まず非同期型の場合、複数のノードが関連する異なるファイルを順に更新した場合、更新の伝播の順序性が保証されない。その為、参照だけを行うノードだけからみても、新旧入り交じった一貫性の無いデータが見えてしまうという本質的な欠点を抱えている。

【 0 0 0 6 】

また複数のノードが同じファイルをほぼ同時に更新（実時間ではかなりずれている場合を含む）すると、各ノードが異なったデータを保持することになり、結果的にファイルが破壊される。

【 0 0 0 7 】

このデータの破壊については、非同期伝播を用いた場合のみならず同期型伝播を採用した場合においても、2つのノードがほぼ同じ時刻に更新した場合にはファイルが破壊されることがある。例えば同じファイルの重複する領域をノードAとノードBの2つのノードがほぼ同じ時刻に更新した場合、ノードAとBが異なるデータを保持する結果になることがある。この場合、その後の処理においては各ノードは自身が保持する互いに異なるデータに基づいて処理を続行することになるので、最終的には矛盾した処理がノードAとノードBで行われることになる。

【 0 0 0 8 】

この為、従来のファイルレプリケーション方式では、スタティックに決まる1つのノードにのみファイルの更新処理を許し、他ノードに対しては参照のみを許すという制約を与えていた。この方式によるものとしては、特開平9-9118

5号「分散コンピューティングシステム」がある。ここで提案されている方式では、自ノードのデータを更新あるいは参照できることを示すWrite トークンと、自ノードのデータを参照のみできることを示すReadトークンを用意し、Write トークン保持ノードが存在する時には、他のノードのいずれもRead/Writeトークンを保持していないように制御し、かつ更新要求を全て同期型で行うことで、同時更新に伴う矛盾を解消している。

【0009】

【発明が解決しようとする課題】

しかし上記公報に開示されている方式では、ファイルの更新が常に同期型で行われるため、応答性の問題を持つ。また、同じファイルを同時にアクセスする複数のノードが存在し、かつその中に1つ以上がファイルの更新を行うものであった場合、自ノードのデータをアクセスする為に必要となるトークンの取り戻し処理をアプリケーションプログラムがI/O要求を発行する度に行わなければならない、オーバーヘッドが非常に大きくなってしてしまう。

【0010】

また、この方式を含め従来のレプリケーション方式では、アクセスは常に自ノードが保持しているデータに対して行われることを前提としているため、新規ノードを系に組込む場合、新たに組込まれたノードは関連する全てのファイルのデータを系内の他のノードから自ノードに一括して取込んだ後でないとデータの一貫性が保証されない。この為、新規ノードは系に加わっても直には業務に移れない、新規参入ノードへのデータ取込み中既存系の更新が止まる、即ち長時間通常業務が停止するという欠点を持っていた。

【0011】

本発明は、最新データを保持しているノードを特定し、そのノードにRead/Write要求を伝播してファイルアクセスを依頼することにより、新規ノード参入時の稼動中業務への影響を最小化することが出来るファイルレプリケーションシステムを提供することを課題とする。

【0012】

また非同期型の伝播においても複数ノードでの同時更新が可能な高速レプリケ

ーションを実現するファイルレプリケーションシステムを提供することを課題とする。

【0013】

更に本発明では、非同期に送られる更新要求のファイルへの反映を、Write 要求のノード内順序性を示す更新番号とそのWrite 要求が前提とする他のノードの更新番号からなる依存ベクタを用いて制御することで、系縮退時でもファイル更新の論理的順序性を保証するファイルレプリケーションシステムを提供することを課題とする。

【0014】

【課題を解決するための手段】

図1は本発明のによるノードの原理図である。

本発明によるノード1は、他のノードとネットワークによって接続され、該他のノードとの共用ファイル6を保持することを前提としており、I/O要求インタセプト部2及びトークン管理手段3を備える。

【0015】

トークン管理手段3は、上記共用ファイル6に対するアクセス要求を管理する。

I/O要求インタセプト手段2は、自ノード内で生じた共用ファイル6に対するアクセス要求に対し、上記トークン管理手段に該共有ファイルへのアクセス許可を求め、許可が得られると共用ファイル6へアクセスする。

【0016】

上記トークン管理手段3は、上記I/O要求インタセプトからのアクセス許可に対し、既に他のノードが上記共用ファイルに対する更新許可を保持する時、該更新許可を保持するノードを上記I/O要求インタセプト手段に通知し、上記I/O要求インタセプト手段は、上記アクセス許可が得られない時、該更新許可を保持するノードに上記共用ファイルへのアクセス処理を依頼する。

【0017】

これにより、各ノード1は共用ファイル6へのアクセスを、最新のデータを保持しているノードのデータに対して行え、また各ノードからは一貫性の有るデー

タが見える。

【 0 0 1 8 】

ノード 1 は、また上記共用ファイル 6 への更新時に更新内容を他の更新との依存関係を示す情報と共に他のノードへ伝播する変更データ通知手段 4 と、上記依存関係を示す情報に基づいて、更新の順序性保証をしつつ上記更新内容を上記共用ファイルに反映させる受信データ処理手段 5 を更に備える構成とすることも出来る。

【 0 0 1 9 】

この構成により、ファイルの更新内容が更新順と前後して到着しても、共用データ 6 は、順序性を保証された更新が行われる。

更にノード 1 は、新規参入時に自ノードの保持する共用ファイル 6 のデータの復元処理を行う系構成管理手段を更に備える構成とすることも出来る。この構成の場合、上記ファイルの復元処理中に、自ノード内で上記共用ファイルに対するアクセス要求が生じた時、上記 I O 要求インタセプト手段 2 は、上記共用ファイル 6 を共用している他のノードにアクセス処理を依頼する。

【 0 0 2 0 】

これの構成により、新規参入したノードは共用ファイルの更新処理の完了を待たずに他の処理に移れる。

【 0 0 2 1 】

【発明の実施の形態】

以下に本発明に於ける一実施形態について図面を参照しながら説明する。

本実施形態のファイルレプリケーションシステムは、複数のノードがネットワークに接続されて系を構成し、系内の各ノードがファイルを共用する構成を前提としている。

【 0 0 2 2 】

まず本実施形態での系の構成について説明する。

図 2 は本実施形態での系及び系の再構成を説明する図である。

本実施形態で系とは、同一のファイル群（以下各系で共用している 1 乃至複数のファイル（群）をオブジェクトグループという）を共用しているノードのグル

ープを指す。例えば図2では、オブジェクトグループa, dを共用しているノードA、C、E及びFで構成される系a、オブジェクトグループbを共用しているノードA、B、及びDで構成される系b、オブジェクトグループcを共用しているノードG、H及びIで構成される系cの3つの系が構成されている。

【0023】

この系内のノードの内、1つのノードが系内の共用ファイルへアクセスするためのRead/Write トークンを管理している。このトークンを管理するノードには、系を構成する際に、予め決められているノードになるか、特定の条件、例えば最小ネットワークアドレスを持つものが動的に選ばれる。

【0024】

また新規のノードが系に加わったり、構成要素となっているノードやネットワークの障害等で系の縮退が生じた時、系の再構成が行われる。例えば、図2の場合系aではノードEの障害によりノードE及びFがネットワークから脱落して残りのノードによる系の再構成が行われている。また系cでは、ノードJがJoinコマンドによって新規に系に加わったことにより系の再構成が行われる。この系の再構成の際には、新規参加ノードの共用ファイルの一貫性 (consistency) 保証のため等価性回復処理が行われる。

【0025】

尚ノードの系からの離脱は、障害などによって生じるものの他、離脱を行うノードが系内の他のノードにメッセージを送信して自律的に行うものが有る。

図3は本発明に於けるノード間の基本動作を説明する図である。

【0026】

図3 (a) は、オブジェクトグループに対してアクセスする際の、ノード間の処理を示す図で、同図では同一の系にノードA～Eの5つのノードあり、そのうちノードAがトークン管理ノードとする。各ノードはユーザプログラムからオブジェクトグループ内のファイルに対するアクセス要求が生じると、ノードAにRead/Writeトークンの獲得要求を発行する。

【0027】

これに対しノードAは、他のノードに既にWrite トークンを渡していなければ

、要求されたトークンを与える。またもし既にWrite を他のノードに渡していれば、トークン獲得失敗通知と共にWrite トークンを保持しているノードを通知する。トークン獲得失敗を通知されたノードは、ファイルへのRead/Write要求をこの通知されたノードに対して依頼し、Write トークンを保持しているノードは、ファイルの順序性を保つようにこれらの要求を処理してゆく。同図の場合ノードB、CがRead要求（参照要求）を、ノードDがWrite 要求（更新要求）を発行した時点でWrite トークンはEが保持しているので、ノードAは各ノードからのトークン獲得要求に対して、獲得失敗と共にノードEがWrite トークンを保持していることを通知する。これに対して各ノードはノードEに対してファイルへのRead/Write要求を発行し、ノードEはファイルの順序性を保持しながらこれらの要求に対しファイルへのRead/Writeを行う。

【 0 0 2 8 】

この様に本発明では、共用ファイルへのアクセス要求が生じたノードに対しWrite トークン保持ノードの通知という形で、その共用ファイルに対する最新のデータを保持するノードが通知される。よって、共用ファイルをアクセスするノードは常に最新のデータに対してアクセスすることが出来る。

【 0 0 2 9 】

また各ノードは、トークンの獲得に失敗してもトークンを獲得できるまで待つことなく処理を続行できる。更に複数のノードによる同一のファイルに対する同時アクセスを可能としている。この為、高い反応性を持つシステムを構築することができる。

【 0 0 3 0 】

またファイルへの変更処理は、他のノードで発生した更新要求に対する処理もWrite トークンを持つ1つのノードが行うので、各ノードからは一貫性のあるデータが見える。

【 0 0 3 1 】

更に同時に生じたアクセス要求を処理する際、それぞれに対するトークンの回収処理を行う必要が無く、オーバーヘッドを小さくすることが出来る。

次に本発明に基づいたシステムに於ける系への新規参入時処理については説明

する。

【 0 0 3 2 】

図 3 (b) は、系に新たに加わったノードの系内の他のノードからとの処理を表す図である。

本発明では、各ノードはデータの最新性を示す情報を保持しており、新規参入ノードはこの情報を比較して、自己が系から離脱している間にデータが更新された時のみ復元処理を行う。また新規ノードはデータの復元処理中に、ユーザプログラムを起動し通常業務に入る。そしてユーザプログラムからファイルへのアクセス要求が発生した場合には、系内の他のノードにRead/Write要求を発行し、ファイルへのアクセス処理を依頼する。図 3 (b) では、系に新規参入したノード D はファイルの復元処理の完了を待たずにユーザプログラムを起動して、ファイルの復元処理中にユーザプログラムからオブジェクトグループ内のファイルへのアクセス要求が生じると、このアクセス要求をWrite トークンを保持してるノード E に依頼している。

【 0 0 3 3 】

この様に本発明では、新規ノードはファイルの復元処理の完了を待たずに、ファイルへのアクセスを行うことが出来るので、系への参入後直ちにユーザプログラムを起動して通常処理を開始することが出来る。

以下に上記基本原理を実現するための一実施形態について図面を参照しながら説明する。

【 0 0 3 4 】

図 4 は本実施形態の系を構成する複数のノードの内の 1 つの構成を示すブロック図である。

システム内の複数のディスク装置上に置かれるオブジェクトグループを共用する各ノード 1 0 は、系構成管理部 1 1、I O 要求インタセプト部 1 2、トークン管理部 1 3、変更データ通知部 1 4 及び受信データ処理部 1 5 が配置される構成となっている。これらの各構成要素は、各ノード内でメモリ上に展開されるプログラムによって実現される。また処理速度を得る為、一部をハードウェアにより

実現する構成としてもよい。また、ノード 10 のローカルディスク装置 18 には、同一の系内で共用している共用ファイル 19 及び系構成の為の定義情報である環境定義・状態情報 20 を記憶している。

【 0 0 3 5 】

尚これらの構成要素の内、I/O 要求インタセプト部 12 はオペレーションシステム (OS) の一部として動作し、ユーザプログラム 17 が発行した入出力命令を受取り、OS 内のファイルシステムにこの入出力命令を伝える役割をする。

【 0 0 3 6 】

尚、本実施形態では、I/O 要求インタセプト部 12 を OS のファイルシステム 16 と分離した構成としているが、ファイルシステム 16 内に含める構成とすることも可能である。また他の構成要素は、OS 内の要素として構成としてもよいし、アプリケーションプログラムとして OS 上に実装する構成としてもよい。

【 0 0 3 7 】

以下、各構成要素について詳細に説明する。

[系構成管理部]

系構成管理部 11 は、ノード起動時や系再構成時における系構成状態の監視、対象ファイルや伝播モードについての設定、ノード障害などに伴う系の縮退や新規ノードの参入等系の状態管理、系再構成時の他ノードとの同期 (同期回復)、新規参入ノードの初期同期 (等価性回復)、ノードの状態の監視及びオペレータとのインタフェース処理を司る部分である。

【 0 0 3 8 】

また系構成管理部 11 は、Join コマンドにより系に加わり Leave コマンドにより系から脱退するまで、後述する系を構成するノードのノード障害監視処理を行う。

【 0 0 3 9 】

システム立ち上げの一環としてファイルレプリケーションシステムを実現するプログラムが起動されると、まず環境定義・状態ファイルが読み込まれ、複数の対象とするオブジェクトグループに属するファイル群、そのオブジェクトグループを配置するノード群、及び更新データの伝播モードについての情報を得る。

【 0 0 4 0 】

この環境定義・状態ファイルは、各オブジェクトグループ毎に構成された系状態テーブルによって構成されている。

図 5 は系状態テーブルの構成例を示す図である。

【 0 0 4 1 】

各系状態テーブルは、オブジェクトグループ毎にそのオブジェクトグループの構成等の情報を記録したテーブルである。各系状態テーブルにはそのテーブルに情報が記憶されているオブジェクトグループを識別するオブジェクトグループ番号、系のバージョン番号、自己が前回整然停止したかどうかを表示する整然停止フラグ、系を構成する各ノードを特定する複数のノード番号とそのノードが前回整然停止したかどうかを示すフラグとからなる複数の配列で構成されるノード定義部、このオブジェクトグループに属する各ファイルを特定するオブジェクトグループ定義部及びこのオブジェクトグループに属するファイルの更新データ伝播モード（同期、半同期、非同期：これらの詳細については後述する）を指定する情報によって構成されている。尚「整然停止」とは、例えば正月休み等でサービスを休止する時に、系内のノードが同期を取って同時にそのオブジェクトグループに対する処理を停止する系からの離脱の仕方を指す。

【 0 0 4 2 】

尚図 5 中の * 部分の情報は、初期値はユーザが設定し、以降系構成管理部 1 1 が必要に応じて変更する情報である。また * が記されていない部分は、ユーザは設定を行わず系構成管理部 1 1 のみが設定、変更する情報であることを示している。

【 0 0 4 3 】

環境定義・状態情報 2 0 は、複数の系状態テーブルからなる構成であり、複数のオブジェクトグループそれぞれに対して設定を行うことが可能である。よって、オブジェクトグループ毎に異なるノード群や更新データの伝播モードを設定することが出来る。例えば、図 2 において、ノード A は、オブジェクトグループ a、b 及び d の 3 つのオブジェクトグループに対する系状態テーブルを持ち、それぞれに異なったノード群（オブジェクトグループ a 及び d にはノード C、D、E

及びF、オブジェクトグループbにはノードB、C及びD)と転送方式(同期, 非同期, 半同期)を設定することが出来る。そして、データの重要度に応じて、例えば、最も重要なオブジェクトグループaには同期モード、重要度の低いcには非同期モード、その中間のオブジェクトグループbには半同期モード等のそれぞれのオブジェクトグループ毎に異なった設定を行うことが出来る。

【0044】

環境設定部は、この環境定義・状態情報20を読み込んで、メモリ上に内部制御表を各オブジェクトグループ毎に展開し、各構成要素にユーザが指定した設定を伝える。

【0045】

この内部制御表は、ユーザが設定したオブジェクトグループの情報を保持するノードのメモリ上に展開されるテーブルで、例えば、図6の様な構成を取る。

図6の内部制御表は、各オブジェクトグループを特定するオブジェクトグループ番号、更新データのデータ伝播モード(同期, 非同期, 半同期)、状態フラグ、オブジェクトグループ定義部、ノード定義部、及び更新伝播送信キューと実反映遅延キューのエントリを示すポインタを記録している。このうちオブジェクトグループ定義部は、系状態テーブルのオブジェクトグループ定義部と同様、そのオブジェクトグループに属するファイル群を特定する先頭ファイルパス名の集合を保持しており、この中に特定されたパス名で始まるファイル群がこのオブジェクトグループに属することを示す。またノード定義部内のノード番号とstatusからなる配列は、このオブジェクトグループを配置するノード群とその状態(動作中, Join中等)を示している。尚更新伝播送信キューと実反映遅延キューについては後述する。

【0046】

また状態フラグは、オブジェクトグループに属するファイルへのアクセス可否や、等価性回復中、系再構成中等の状態を表示するフラグの集合で、図4に示した各構成要素はこの状態フラグの対応ビットの1/0を切換えることよりこれらの状態を表示して他の構成要素に通知する。尚初期状態では、既に他のノードが系を作り、ファイルが更新されている可能性があるのでオブジェクトグループに

属するファイルは全てアクセス不可の状態として設定される。

【 0 0 4 7 】

初期処理が完了すると、系構成管理部 1 1 はオペレータからオブジェクトグループに対する操作指令が投入されるのを待つ。

1) Joinコマンド投入

オペレータはオブジェクトグループに対する活性化を指示する場合、Joinコマンドを投入する。

【 0 0 4 8 】

このJoinコマンドが投入されると、系構成管理部 1 1 は、他のノードとメッセージをやり取りしてJoinコマンドと共に指定されたオブジェクトグループに対する系に加わる。またJoinコマンドに単独での系生成を許可することを示すsingle指定がされていた場合、もしこのオブジェクトグループに対して系が構成されていないければ新たな系を生成する。

【 0 0 4 9 】

図 7 は、Joinコマンド投入時の系構成管理部 1 1 による処理を示すフローチャートである。

Joinコマンドが投入されると、系構成管理部 1 1 は、まずJoinコマンドと共に指定されたオブジェクトグループを共用している他のノードに順にメッセージを送り（ステップ S 1 1）、返答を各ノードから受信する（ステップ S 1 2）。

【 0 0 5 0 】

各ノードからの返答から、対象としているオブジェクトグループに対して既に系を作っているものでないかどうかを調べ、その結果既に系を作成しているという返答が他ノードからあれば（ステップ S 1 3、YES）、そのノードにJOIN要求を送り既存系への参入処理を依頼する（ステップ S 1 4）。

【 0 0 5 1 】

この参入依頼に対し、ノードから参入失敗を通知する応答がされた場合（ステップ S 1 5、YES）、Join失敗をオペレータに通知し（ステップ S 1 6）、処理を終了する。また参入失敗の通知がなければ（ステップ S 1 5、NO）、後述の参入処理（ステップ S 1 7）を行った後、オペレータに成功応答を返す（ステ

ップS18)。

【0052】

またステップS12の各ノードからの応答から、そのオブジェクトグループに対して未だ系を作っているノードがいないと判断され(ステップS13、NO)、かつJoinコマンドのオプションでsingleが指定されていた場合(ステップS19、YES)、このノードは自身のみで系を作る。

【0053】

この際、系構成管理部11は、まず、系状態テーブル内の情報を調べる。その結果、系状態テーブル中の整然停止フラグに整然停止が表示され自身が認識している最終の系状態が整然停止であると判断される時(ステップS20、YES)、一定時間受信待機し(ステップS21)、前回整然停止した時に共に系を構成していた他のノードが新規系への参入を依頼してくるのを待つ。そしてJOIN要求により系への参入を依頼してきたノードに対し、順次後述する図9のJOIN要求受け付け処理を行い、自己の系のバージョン番号を送信する。

【0054】

この結果、全てのノードからREADY要求が到着したら(ステップS22、YES)、READY要求に対する応答として、COMPLETE応答を全ノードに返す(ステップS23)。また全てのノードからREADY要求が到着しなければ(ステップS22、NO)、READY要求に対する応答としてCONT応答をノードに返し(ステップS24)、更にREADY要求の到着を待つ。

【0055】

ステップS23でREADY要求に対する応答を送信した後、あるいはステップS20で系状態テーブル内の整然停止フラグが前回の停止が整然停止でないことを示していた場合(ステップS20、NO)、環境定義・状態情報20の対応する系状態テーブルに記録されている系のバージョン番号をインクリメント(+1)して更新する(ステップS25)。そして、内部制御表の状態フラグをアクセス可能表示に変更して(ステップS26)、IO要求インタセプト部12に対応するオブジェクトグループへのアクセスが可能となったことを知らせる。そしてJoinコマンドに対する応答として処理完了をオペレータに通知して(ステップ

S 2 7) 処理を終了する。

【 0 0 5 6 】

またステップ S 1 9 で、Join コマンドのオプションとして single 指定がされていなかった場合には（ステップ S 1 9、NO）、Join コマンドに対する応答としてオペレータにエラーを通知し（ステップ S 2 8）、処理を終了する。

2) 参入処理

図 8 は、図 7 のステップ S 1 7 の系構成管理部 1 1 の動作を示すフローチャートである。

【 0 0 5 7 】

J O I N 要求による参入依頼に対し、参入失敗でなければ依頼先のノードから応答として系のバージョン番号が送信されてくる。この時系構成管理部 1 1 は、まず現在系を構成するノード情報から内部制御表の依頼元ノードに対応する status を Join 中表示に更新し（ステップ S 3 1）、次に応答で通知された既存系が保持しているバージョン番号と参入しようとしている自ノードが保持しているバージョン番号を比較する（ステップ S 3 2）。その結果、2つのバージョン番号が異なる場合には、自ノードが系から脱落している間にオブジェクトグループ内のファイルに対し変更が加えられた可能性があることを示しているので、整然停止表示をリセットし（ステップ S 4 1）、等価性回復処理を起動する（ステップ S 4 2）。また2つの系のバージョン番号が一致していても系状態テーブルの整然停止フラグが非整然停止を表示していた場合には（ステップ S 3 2、一致：ステップ S 3 3、NO）、自己のファイルは最新のデータのものでないので、やはりステップ S 4 2 の等価性回復処理を起動する。そしてステップ S 4 2 の等価性回復処理の起動後は、処理の終了を待たずに応答値として送信されてきた系のバージョン番号を系状態テーブルに設定し（ステップ S 4 3）た後、内部制御表の状態フラグをオブジェクトグループに対するアクセスが可能の表示に変更し（ステップ S 4 0）、処理を終了する。

【 0 0 5 8 】

また送信されてきた系のバージョン番号と系状態テーブル内に記憶されている系のバージョン番号が一致しており（ステップ S 3 2、一致）、かつ系状態テ

ブルの整然停止フラグが整然停止を表示しているなら（ステップ S 3 3、YES）、自ノードが保持しているオブジェクトグループのファイルは最新のデータものなのでファイルの更新は必要ない。よって後述するステップ S 4 2 の等価性回復処理は行われず、ステップ S 3 4 として系のバージョン番号を更新後、定期的に READY 要求を送り（ステップ S 3 5）、全ノードの参入が完了するのを待合わせる。

【0059】

その結果 READY 要求に対する応答が CONT 応答であれば（ステップ S 3 6、CONT）、一定時間後に READY 要求を再送し（ステップ S 3 7）、同じ処理を繰り返す。また READY 要求の応答が COMPLETE 応答であれば（ステップ S 3 6、COMPLETE）、前回整然停止した時のノードが全て依頼元に READY 要求を行ったことになるので、応答で返される系を構成しているアクティブノードについての情報から、内部制御表のノード定義部の各ノードの status を動作中表示に変更する（ステップ S 3 8）。

【0060】

この後、系状態テーブルの整然停止表示をリセットし（ステップ S 3 9）、内部制御表の状態フラグをオブジェクトグループがアクセス可能表示に変更し（ステップ S 4 0）、処理を終了する。

3) JOIN 要求受付処理

図 9 は、JOIN 要求受付処理時の系構成管理部 1 1 の動作を示すフローチャートである。

【0061】

この JOIN 要求受付処理は、図 7 のステップ S 1 4 の新規参入依頼時に発行された JOIN 要求や、ステップ S 2 1 の受信待機時に受け付けた JOIN 要求に対する処理を示したものである。

【0062】

JOIN 要求を行ったノードから受け付けたノードの系構成管理部 1 1 は、JOIN 要求と共に通知された依頼ノードの系のバージョン番号と系状態テーブル内の自身の系のバージョン番号とを比較する（ステップ S 5 1）。その結果両方

のバージョン番号が一致しており（ステップ S 5 1、一致）、また整然停止フラグを参照して整然停止後の整然立ち上げ中なら（ステップ S 5 2、YES）、ステップ S 5 3 として JOIN 要求に対する応答して自己の現在のバージョン番号を返答して処理を終了する。

【 0 0 6 3 】

JOIN 要求と共に通知された情報から、2つの系のバージョン番号が一致しなかったり（ステップ S 5 1、不一致）、一致しても整然停止後の系への参加でないのならば（ステップ S 5 2、NO）、次に内部制御表のノード定義部を参照し、既にJoin中のノードが存在するかどうかを調べる（ステップ S 5 4）。その結果、既にJoin中のノードが存在していれば、応答として失敗を通知して（ステップ S 5 9）、処理を終了する。またJoin中のノードが他に存在しなければ（ステップ S 5 4、NO）、この JOIN 要求により参入してきたノードに対応する内部制御表内のstatusを稼動中（アクティブ）、JOIN 中（新規参入処理中）の表示に設定し（ステップ S 5 5）た後、他のアクティブな全ノードにJoin通知を送る（ステップ S 5 6）。そしてこのJoin通知に対する応答が全て返ってきた後に（ステップ S 5 7、YES）、系のバージョン番号を更新し（ステップ S 5 8）、JOIN 要求に対する応答として現在の系のバージョン番号を返答して（ステップ S 5 3）、処理を終了する。

4) Join通知

図 1 0 は、図 9 のステップ S 5 6 で送信されたJoin通知を受取ったアクティブなノードの系構成管理部 1 1 が行う処理を示すフローチャートである。

【 0 0 6 4 】

Join通知を受信すると、系構成管理部 1 1 は、ステップ S 6 1 として内部制御表の、Join通知により通知された参入依頼をしているノードに対応するstatusを稼動中、Join中表示に設定する。そしてステップ S 6 2 として、Join通知に対する応答後、系状態テーブル内の系のバージョン番号を更新して（ステップ S 6 3）処理を終了する。

5) 等価性回復処理

図 1 1 は、図 8 のステップ S 4 2 で起動される等価性回復処理の系構成管理部

1 1 の動作を示すフローチャートである。

【 0 0 6 5 】

等価性回復処理は、新規参入ノードが系から離脱している際に古くなった自己のファイル内のデータを復元する為の処理である。

等価性回復処理が起動されると、まず系構成管理部 1 1 は内部制御表のノード定義部を参照して、系内のアクティブなノードの 1 つからオブジェクトグループ内の全ファイルのファイル名を取得する（ステップ S 7 1）。

【 0 0 6 6 】

次にステップ S 7 2 として内部制御表の状態フラグを等価性回復中表示に設定した後、ステップ S 7 3 系内のアクティブなノードにステップ S 7 1 で得たファイル名を指定して転送要求を行う。この転送を等価性回復転送と呼ぶ。

【 0 0 6 7 】

このファイル転送に対する応答がエラーであったならば、転送要求先を他のアクティブなノードに変更して再度ファイル転送要求を行う（ステップ S 7 5）。

ファイル転送要求に対して、要求先のノードから、正常応答を得たら（ステップ S 7 4、正常）、ステップ S 7 5 として転送ファイルを受信し、これを受信データ処理部 1 5 に自身のファイルにデータの反映を依頼する（ステップ S 7 7）。この時通常のファイル更新に伴う更新データの伝播と、等価性回復処理での転送データの順序性は変更データ通知部 1 4 及び受信データ処理部 1 5 を介して保証されるので、等価性回復処理中にファイルを更新しても更新結果が失われることはない。

【 0 0 6 8 】

転送ファイルの受信及び自ファイルへの反映をステップ S 7 1 で得た全てのファイルに対して行い（ステップ S 7 8、NO）、全ファイルへの処理が完了したならば（ステップ S 7 8、YES）、ステップ S 7 9 として全アクティブノードに等価性回復処理の完了を通知し、全アクティブノードからの応答を待った後（ステップ S 8 0）、内部制御表上の等価性回復処理中をリセットし（ステップ S 8 1）、処理を終了する。尚ステップ S 7 3 ～ 7 8 の等価性回復転送によるファイル転送は 1 つのノードに全ファイルの転送を要求してもよいし、複数のノード

に分散して要求してもよい。

6) 等価性回復転送

図 1 2 は、等価性回復処理を行っているノードから、図 1 1 のステップ S 7 3 で送信される等価性回復転送要求を受信したノードの系構成管理部 1 1 が行う処理を示すフローチャートである。

【 0 0 6 9 】

等価性回復転送を要求されたノードは、ステップ S 9 1 としてまずトークン管理ノードに Write トークンの獲得を要求する。その結果、Write トークンを獲得できなければ（ステップ S 9 2、N O）、要求先のノードにエラー応答をして（ステップ S 9 3）処理を終了する。

【 0 0 7 0 】

また Write トークンを獲得できれば（ステップ S 9 2、Y E S）、ステップ S 9 3 として要求元のノードに正常を応答した後、ステップ S 9 5 として要求されたファイルデータを変更データ通知部 1 4 を介して順次要求元のノードに転送し、処理を終了する。

7) 等価性回復完了メッセージ

図 1 3 は、等価性回復処理が完了した最新のデータにファイルの復元がなされたノードが、図 1 1 のステップ S 7 9 で送信した等価性回復完了メッセージを受信した系内のアクティブなノードが行う処理を示すフローチャートである。

【 0 0 7 1 】

等価性回復完了メッセージを受信したノードは、ステップ S 9 6 として内部制御表内の送信もとノードに対応する status に表示されている J O I N 中の表示をリセット後、ステップ S 9 7 としてメッセージの送信元ノードに応答を返して処理を終了する。

【 0 0 7 2 】

この図 1 3 の処理により、新規参入してきたノードは系内の他のアクティブノードから系への参入処理が完了したものとみなされる。

8) Join再試行メッセージ

Join中に系の再構成が発生すると、このJoin再試行メッセージがJoin中のノード

ドに送られる。この要求を受けた系構成管理部 1 1 は、系への新規参入処理を最初からやり直す。

9) 停止処理

ノードを停止させる場合、オペレータは系からの離脱を指示する leave コマンドを投入して系から離脱する。尚ここでの停止とは、ノードが系から離脱することを示しており、ノードが複数の系に属している場合、メンテナンス等でノードを完全に止めるためには各系に対して leave コマンドを投入して全ての系から離脱しなければならない。

【0073】

ノードの停止を leave コマンドでオペレータから通知されると、系構成管理部 1 1 は以下の処理を行う。

a) 整然停止

整然停止は、系を構成している全ノードが同期して一斉に停止し系そのものが停止するもので、正月休みやシステムの再構築等の場合にシステム全体を休止させるために行われる。オペレータは整然停止を行う場合、オプションで all を指定した leave コマンドを投入する。

b) 非整然停止

非整然停止は、そのノードのみを停止させるものであり、非整然停止したノードのみ系から離脱し、他のノードによって系は存続する。オペレータは非整然停止を行う場合、オプションで all を指定しないで leave コマンドを投入する。

【0074】

図 1 4 は、オペレータが leave コマンドを投入して、ノードの停止を指示した時の系構成管理部 1 1 の処理を示すフローチャートである。

leave コマンドが投入されると、系構成管理部 1 1 は、まずステップ S 1 0 1 として内部制御表の状態フラグをアクセス不可表示に変更し、図 4 の他の構成要素に（具体的には I O 要求インタセプト部 1 2 に）対応するオブジェクトグループに属するファイルへのアクセスを禁止する。

【0075】

次に系構成管理部 1 1 は、ステップ S 1 0 2 として変更データ通知部 1 4 に S

YNC 要求を行い、キューに保持され遅延している更新要求の全ノードへの反映を依頼する。

【0076】

全ノードへの変更データの反映が完了し、変更データ通知部 14 から完了が通知されると（ステップ S103、YES）、leave コマンドに all 指定が無い場合には（ステップ S104、NO）、非整然停止なので処理を終了する。

【0077】

またステップ S104 で all 指定がある場合には、整然停止を行うため、ステップ S105 として整然停止開始メッセージを系内の全ノードへ一定時間送信し、整然停止開始メッセージに対する応答が全ノードから返信されるのを待つ（ステップ S106）。そして全ノードから応答があると（ステップ S106、NO）、整然停止を行ったオブジェクトグループに対応する系状態テーブル内の整然停止フラグを整然停止にセットして（ステップ S107）、処理を終了する。

10) ノード障害認識

障害等による他のノードの離脱は、例えば分散システムで一般的に行われている自己の存在を他ノードに通知するメッセージ（I'm alive メッセージ）を送信し合うグループコミュニケーションシステムにおいて、メッセージが途絶えたり、応答が返ってこない等の場合に、系内の他のノードによって認識される。系内の他ノードの離脱を認識したノードは、系の再構成を系内の他のアクティブなノードに要求する。

【0078】

図 15 は、系内の他ノードの離脱を認識したノードの系構成管理部 11 の処理を示すフローチャートである。

現在系を構成しているノードの障害を認識すると、系構成管理部 11 はまずステップ S111 として、内部制御表の状態フラグを系再構成中を表示するよう設定し、変更データ通知部 14 にメッセージを他のノードに送るのを一時抑止させる。

【0079】

次に系構成管理部 11 は、ステップ S112 として、系の再構成要求メッセー

ジを系内の全アクティブノードに送信して他のノードの系構成管理部 1 1 とやり取りし、系の再構成の合意を得る。この時、もしJoin中のノードを除く過半数のノードから合意が取れなければ（ステップ S 1 1 3、N O）、状態フラグをアクセス禁止の表示にセットして（ステップ S 1 1 4）、対応するオブジェクトグループ内のファイルへのアクセスを禁止した後、ステップ S 1 1 1 でセットした系再構成中の表示をリセットして（ステップ S 1 1 5）、処理を終了する。

【 0 0 8 0 】

またステップ S 1 1 3 で、系の再構成要求に対してJoin中のノードを除く過半数のノードから合意が取れると（ステップ S 1 1 3、Y E S）、系状態テーブル内の系のバージョン番号を更新し（ステップ S 1 1 6）、ノード定義部の各ノードのstatusを変更して合意の取れた過半数のノードを新しいアクティブなノードとして内部制御表に設定して（ステップ S 1 1 7）、最新の系状態を表すように更新する。

【 0 0 8 1 】

この後、変更データ通知部 1 4 に R E S E T 要求を送り（ステップ S 1 1 8）、応答を待つ（ステップ S 1 1 9）。変更データ通知部 1 4 から応答があったら（ステップ S 1 1 9、Y E S）、更新伝播送信キュー内の変更内容の他ノードへの伝播完了を通知する R E S E T C O M P をアクティブな全ノードの系構成管理部 1 1 に送り、全ノードから R E S E T C O M P に送られてくるのを待つ（ステップ S 1 2 1）。

【 0 0 8 2 】

全ノードから R E S E T C O M P が送られてきたら（ステップ S 1 2 1、Y E S）、伝播中であったファイルの更新要求が全て自ノードに到着したことになるので、ステップ S 1 2 2 として受信データ処理部 1 5 に R E S E T 要求を送り、系から切り離されたノードに関する送信、受信の後始末を依頼し処理完了通知を待つ（ステップ S 1 2 3）。

【 0 0 8 3 】

受信データ処理部 1 5 から処理の完了が通知されると（ステップ S 1 2 3、Y）、ステップ S 1 1 1 でセットした系再構成中の表示をリセットして（ステップ

S 1 2 4)、処理を終了し、通常処理を再開させる。

【 0 0 8 4 】

尚Join中のノードには、Join再試行要求を送り、最初から系への新規参入処理をやり直させる。

[I O 要求インタセプト部]

I O 要求インタセプト部 1 2 は、ユーザプログラム 1 7 が発行したファイルへのアクセス要求を受取り、O S 内のファイルシステムにこのアクセス要求を伝える部分で、ユーザプログラム 1 7 がファイルに対する入出力要求を発行すると、I O 要求インタセプト部 1 2 に制御が渡る。

【 0 0 8 5 】

I O 要求インタセプト部 1 2 は要求されたファイルの名前が全ての内部制御表に設定されているいずれのパス内にも属していないなら、直ちにO S のファイルシステムに制御を渡す。そしてファイルシステムから戻された応答をユーザプログラム 1 7 に返す。

【 0 0 8 6 】

またもしそのファイルが、複数の内部制御表の内のオブジェクトグループ定義部内に定義されているいずれかのパスに属するものであるならば、要求されたファイルへのアクセス要求がオブジェクトグループに属するものと見なし、以下の処理を行う。

1) アクセス不可表示が内部制御表にある場合

オブジェクトグループへのアクセスは禁止されているので、ユーザプログラム 1 7 にエラーを応答する。

2) 等価性回復中の場合

稼動中の他ノードに F O R C E 指定のRead要求若しくはWrite 要求を送り、ファイルへのアクセスを依頼する。系内の他のノード (Join中を除く) は、最新データのファイルを保持しているので、Read/Write要求に対して応答データを送信してきた場合には、このデータは一貫性が保証されているものなのでこれをユーザプログラム 1 7 に返す。またRead/Write要求に対して失敗を応答されたら、別の稼動中ノードに対して同様の処理を繰り返す。

3) 等価性回復中でない場合

a) Write 系要求

要求されたファイルのWrite トークン獲得をトークン管理部 1 3 に依頼する。トークン管理部 1 3 から獲得成功を応答された場合、OS のファイルシステムを呼び、自身のファイルに対しデータの更新処理を行った後、変更内容を変更データ通知部 1 4 に渡して他ノードへの反映を行う。

【0 0 8 7】

トークン管理部 1 3 からWrite トークン獲得失敗を応答されたら、トークン管理ノードから応答時に通知されたWrite トークン保持ノードにWrite 要求を送り処理を依頼する。またWrite トークン保持ノードからWrite 要求に対して処理失敗（トークン変化）を応答されたら、トークン獲得からやり直す。

【0 0 8 8】

尚自ノードのファイルを更新する際の待合わせ処理や、他ノードの受信データ処理部 1 5 に送るWrite 要求に付加するデータなど、I O 要求インタセプト部 1 2 で行われる順序性保証処理は後述する。

b) Read 系要求

要求されたファイルのRead トークン獲得をトークン管理部 1 3 に依頼する。トークン管理部 1 3 から獲得成功を通知されたら、OS のファイルシステムを介し自ノードのファイルからデータを読み、ユーザプログラム 1 7 に応答する。

【0 0 8 9】

トークン管理部 1 3 からRead トークン獲得失敗を応答されたら、応答時に通知されたノード（Write トークン保持ノード）にRead 要求を送る。成功応答が要求先のノードからあれば、渡されたデータをユーザプログラム 1 7 に返す。また失敗（トークン変化）を応答されたらトークン獲得からやり直す。

【0 0 9 0】

尚他ノードで行われた更新の待合わせなど、順序性保証に伴う処理は後述する。

尚Read/Write トークンの獲得／解放はユーザプログラム 1 7 からのRead/Write 要求発行単位で行う構成にする他、オーバーヘッドを減らすためファイルのOpen/C

lose単位に行う構成にしても良い。この場合、ユーザプログラムがファイルをオープンした時点で上記トークン処理が行われ、クローズが発行されるまでトークンが保持される。またトーク獲得不可をオープン時に通知された場合、以降のI/O要求はWrite トークンを保持しているノードに転送される。

【0091】

また、トークン解放を自発的に行うのではなく、ファイル処理が完了するとトークンを必要としていないことを表示しておき、他ノードがトークンを必要とするタイミングまで解放を遅らせる構成とすることも出来る。尚、Write 時及びRead時には後述する順序性保証処理も行われる。

【0092】

図16は、I/O要求インタセプト部による処理を示すフローチャートである。

ユーザプログラム17からファイルへのアクセス要求が発行されるとI/O要求インタセプト部12はまず内部制御表を参照し、要求されたファイルのファイル名とオブジェクトグループ定義部内のパス名を比較する(ステップS131)。その結果一致しなければ(ステップS131、不一致)、要求されたファイルはオブジェクトグループに属していないので、OSのファイルシステムの制御を渡し(ステップS132)、ファイルへの処理を依頼する。そしてファイルシステムからの応答をユーザプログラムに返して(ステップS133)、処理を終了する。

【0093】

ステップS131でファイル名が内部制御表内のいずれかのパスに属するものであるのならば(ステップS131、一致)、そのファイルはオブジェクトグループに属するものである所以对応する内部制御表の状態フラグを調べる。その結果アクセス不可が表示されていれば(ステップS135、YES)、ステップS134としてユーザプログラム17にエラー応答を行い処理を終了する。

【0094】

また状態フラグに等価性回復中表示がされていた場合には(ステップS136、YES)、ステップS150として稼働中の他ノードにオプションでFORCE指定をしたRead/Write要求を送り、応答を待つ(ステップS151)。その結

果失敗を応答されたら（ステップ S 1 5 2、失敗）、ステップ S 1 5 3 として別のアクティブなノードに F O R C E 指定の Read/Write 要求を送り、応答を待つ。また Read/Write 要求を送ったノードから成功応答があると（ステップ S 1 5 2、成功）、ステップ S 1 5 4 として応答データをユーザプログラム 1 7 に応答して処理を終了する。

【 0 0 9 5 】

状態フラグに、アクセス不可と等価性回復中のいずれもが表示されていない時（ステップ S 1 3 5 及び 1 3 6、N O）、I O 要求インタセプト部 1 2 は、アクセス要求が Read 要求であった時（ステップ S 1 3 7、R e a d）、ステップ S 1 4 4 として要求されたファイルの Read トークン獲得をトークン管理部 1 3 に依頼する。

【 0 0 9 6 】

その結果トークン管理部 1 3 からトークンの獲得成功を通知されたら（ステップ S 1 4 5、Y E S）、O S ファイルシステムを介し自ノードのファイルからデータを読みだし（ステップ S 1 4 6）、これをユーザプログラムへ返答（ステップ S 1 4 7）後、トークンを自発的に解放する構成の場合トークン管理部 1 3 にトークンを解放を依頼してから処理を終了する。またトークン管理部 1 3 から Read トークンの獲得失敗を応答された場合には（ステップ S 1 4 5、N O）、ステップ S 1 4 8 として失敗と共に通知された Write トークン保持ノードに Read 要求を送信し、応答を待つ。その結果 Write トークン保持ノードから、成功を通知された場合には（ステップ S 1 4 9、成功）、渡されたデータをユーザプログラム 1 7 に応答後（ステップ S 1 4 7）、処理を終了する。また Write トークン保持ノードから、失敗を通知された場合には（ステップ S 1 4 9、失敗）、ステップ S 1 4 4 の Read トークン獲得依頼から処理をやり直す。尚ステップ S 1 4 6 のデータ読みだし時に読み出し対象となっているファイルの転送モードが非同期若しくは半同期モードである場合、実反映遅延キューを参照して最新のデータがキューイングされていればそちらを読み出す。この点については順序性保証の項で詳細に説明する。

【 0 0 9 7 】

またステップ S 1 3 7 でアクセス要求が Write 要求であった場合には（ステップ S 1 3 7、Write）、ステップ S 1 3 8 として要求されたファイルの Write トークン獲得をトークン管理部 1 3 に依頼する。

【0098】

その結果トークン管理部 1 3 からトークン獲得成功を応答されたならば（ステップ S 1 3 9、YES）、ステップ S 1 4 0 として OS のファイルシステムを呼び自身のファイルに対する Write 処理を依頼し、ステップ S 1 4 1 としてデータの変更内容を変更データ通知部 1 4 に渡して他ノードへの反映を依頼した後、トークン解放を自発的に行う構成の場合トークン管理部 1 3 にトークンの解放を依頼して処理を終了する。またトークン管理部 1 3 から Write トークン獲得失敗を応答された場合には（ステップ S 1 3 9、NO）、ステップ S 1 4 2 として応答時に通知された Write トークン保持ノードに Write 要求を送り、応答を待つ。その結果失敗を応答されたならば（ステップ S 1 4 3、失敗）、ステップ S 1 3 8 の Write トークン獲得からやり直す。また成功を応答されたならば（ステップ S 1 4 3、成功）、後述する順序性保証処理による更新内容の自ファイルへの反映を考慮しつつ処理を終了する。尚ステップ S 1 4 0 で、ファイルに対する処理を依頼する際、対象ファイルの伝播モードが非同期若しくは半同期モードである場合、変更内容は実反映遅延キューにキューイングして順序性保証を考慮した処理が行われる。この点については順序性保証の項で詳細に説明する。

〔トークン管理部〕

トークン管理部 1 3 は、ファイルアクセス権限を管理する部分で、系を構成する全ノードが同じ情報を保持するように制御を行う。尚実装を簡単にする為、系を構成するいずれか 1 つのノード（例えばネットワークアドレスが一番小さいノード）をトークン管理ノードとし、トークン管理ノードのトークン管理部 1 3 をサーバとして系全体の全トークン状態を保持、管理する構成とし、他のノードのトークン管理部 1 3 は、クライアントとして自ノードが保持しているトークンのみを管理する構成とするのが一般的である。

【0099】

トークン管理ノードのトークン管理部 1 3 は、メモリ上にトークン制御表を構

成し、このトークン制御表により系内に存在する全ノードを管理する。

図17は、トークン制御表の構成例を示す図である。

【0100】

同図では、トークン制御表はリストデータ構造を取っており、各オブジェクトグループに属するファイル毎に1つ対応するトークン制御表が生成される。トークン制御表はトークンがオブジェクトグループ内のどのファイルに対してのものであるかを示すファイル識別子、トークンの種類 (Read/Write) を示すトークン状態、トークンを保持しているノードを指定する保持ノード番号及び次の制御表の一を示すポインタが記憶されている。このうちトークン識別子にはトークン管理部13が対応する制御表を検索するためのタグとなるもので、対応するファイルのファイル名等が用いられる。リストの検索を速くするためにファイル識別子にハッシュ関数を適用し、得られた値が同じものが一つのキューを構成するように構成される。

【0101】

トークン管理ノードのトークン管理部13は、自ノードのIO要求インタセプト部12や他ノードのトークン管理部13からトークンに対する処理要求があるところのトークン制御表を検索し、要求されたファイルのトークンの状況を調べる。またトークンを生成したり解放する時は、新たなトークン制御表をリストデータに加えたり、対応するトークン制御表をリストデータから削除する。

【0102】

また系を再構成した場合には、各ノードが保持する最終のトークン保持情報から系全体のトークン状態を復元する。

図18はトークン管理ノードのトークン管理部13の処理を示すフローチャートである。

【0103】

トークン管理部13は、他ノードのトークン管理部13や自ノードのIO要求インタセプト部12からトークンに対する処理要求を受取ると、以下の様に処理する。

【0104】

トークン管理部 1 3 は、他ノードのトークン管理部 1 3 や自ノードの I O 要求インタセプト部 1 2 から処理要求を受取とるとまず、要求内容を判断する（ステップ S 1 5 1）。その結果、Write トークン獲得要求であるのならば、ステップ S 1 5 2 として図 1 9 の Write トークン獲得要求処理を行い、Read トークン獲得要求であるのならば、ステップ S 1 5 3 として図 2 0 の Read トークン獲得要求処理を行い、トークン解放要求若しくはトークン回収要求であるのならば、ステップ S 1 5 4 として、トークン解放／回収要求処理を行ったのち、処理を終了する。

【0 1 0 5】

図 1 9 は、図 1 8 のステップ S 1 5 2 の Write トークン獲得要求処理時のトークン管理部 1 3 の処理動作を示すフローチャートである。

トークン管理部 1 3 は、Write トークン獲得要求処理ではまずトークン制御表を参照して、Write トークン獲得要求を行っているノードが Write トークンを保持しているかどうかを調べる（ステップ S 1 6 1）。その結果、保持していた場合（ステップ S 1 6 1、Y E S）、トークン獲得成功を要求元ノードに応答し（ステップ S 1 6 8）、処理を終了する。また要求元のノードが Write トークンを保持していない場合（ステップ S 1 6 1、N O）、次に要求元以外のノードが要求されているファイルへの Write トークンを保持しているかどうかを判断する。その結果 Write トークンを保持しているノードがある場合（ステップ S 1 6 2、Y E S）、Write トークン獲得不可を Write トークンを保持しているノードのノード番号と共に応答し（ステップ S 1 6 3）、処理を終了する。

【0 1 0 6】

また Write トークンを保持しているノードが存在しない場合には（ステップ S 1 6 2、N O）、他ノードが要求されているファイルの Read トークンを保持しているか判断する。その結果 Read トークンを保持しているノードが存在しなければ（ステップ S 1 6 4、N O）、トークン制御表を操作して要求元ノードに Write トークンを渡し（ステップ S 1 6 8）、トークン獲得成功を要求元ノードに応答して処理を終了する。また Read トークンを保持しているノードが存在すれば（ステップ S 1 6 5、Y E S）、ステップ S 1 6 6 として Read トークンを保持して

全てのノードにトークン回収を指示し、全Readトークン保持ノードから回収完了を通知されるのを待ち（ステップS 1 6 6、NO）、全てのReadトークンの回収が完了した後（ステップS 1 6 6、YES）、要求元ノードにWrite トークンを渡し（ステップS 1 6 7）、トークン獲得成功を要求元ノードに応答して（ステップS 1 6 8）処理を終了する。

【0107】

図20は、図18のステップS 1 5 3のReadトークン獲得要求処理時のトークン管理部13の処理動作を示すフローチャートである。

Readトークン獲得要求処理では、トークン管理部13はまずトークン制御表を参照して、Readトークン獲得要求を行っているノードが、Readトークン若しくはWrite トークンを保持しているかどうか調べる（ステップS 1 7 1）。その結果どちらかのトークンを保持していた場合（ステップS 1 7 1、YES）、トークン獲得成功を要求元ノードに応答し（ステップS 1 7 5）、処理を終了する。また要求元のノードがReadトークン、Write トークン共に保持していない場合（ステップS 1 7 1、NO）、次に要求元以外のノードが要求されているファイルへのWrite トークンを保持しているかどうかを判断する。その結果Write トークンを保持しているノードがある場合（ステップS 1 7 2、YES）、Readトークン獲得不可をWrite トークンを保持しているノードのノード番号と共に応答し（ステップS 1 7 3）、処理を終了する。

【0108】

またWriteトークンを保持しているノードが存在しない場合には（ステップS 1 7 2、NO）、トークン制御表を操作して要求元ノードにReadトークンを渡し（ステップS 1 7 3）、ステップS 1 7 4としてトークン獲得成功を要求元ノードに応答して処理を終了する。

【0109】

図21は、図18のステップS 1 5 4トークン解放／回収要求処理時のトークン管理部13の処理動作を示すフローチャートである。

トークン解放要求は、トークンが不要となったノードが行うもので、系内の全ノードに更新データの伝播が完了した時等に発行される。尚不要になったトーク

ンの解放を自発的に行わない構成の場合、トークン解放要求を受けた時点ではトークン保持ノード側のトークン管理部 1 3 はトークン返却可能を表示しておくのみで復帰する。この場合、Write トークン獲得要求処理やReadトークン獲得要求処理では、トークン管理ノードのトークン管理部 1 3 はWrite トークン保持ノードに対してもトークン回収を指示する。そして、トークンを保持しているノードからトークン回収完了を通知された場合にはトークンを獲得できたとして処理を行い、回収不可を通知された場合にWrite トークンを保持しているノードが存在すると見なして処理を行う。

【0 1 1 0】

またトークン回収要求は、Write トークン獲得要求処理時にトークン管理ノードのトークン管理部 1 3 がRead/Writeトークンを保持しているノードに対して発行する要求である。Write トークンに対する回収要求は、トークン保持ノードが不要になったトークンを自発的に返却しない構成の場合にのみ発行される。

【0 1 1 1】

トークン解放要求若しくはトークン回収要求を受けたトークン管理部 1 3 は、ステップ S 1 8 1 として指定されたトークンを直ちに解放し、解放に成功したことをトークン管理ノードのトークン管理部 1 3 に応答して（ステップ S 1 8 2）処理を終了する。

【0 1 1 2】

図 2 2 は、不要になったトークンを自発的に返却しない構成の場合に発行されるWrite トークン回収要求を受けたWrite トークン保持ノードが行う処理を示すフローチャートである。

【0 1 1 3】

Write トークン回収要求を受けると、Write トークンを解放できる状態にあるかどうか判断される（ステップ S 1 9 1）。その結果、該当ファイルへの書込み処理が完了しておらず、Write トークンを解放できない状態にある場合（ステップ S 1 9 1、NO）、Write トークン回収要求を送信してきたトークン管理ノードのトークン管理部 1 3 にWrite トークン解放失敗を応答し（ステップ S 1 9 6）、処理を終了する。

【 0 1 1 4 】

またWrite トークンを解放可能な状態である時は（ステップS 1 9 1、YES）、まずステップS 1 9 2としてFSYNC指定で変更データ通知部14を呼び、更新伝播送信キューにキューイングされている自己が行ったファイルの変更や他ノードから依頼されているファイルの変更内容を全て、系内の全ノードへの伝播を依頼し、完了応答を待つ（ステップS 1 9 3、NO）。

【 0 1 1 5 】

そして全ノードから応答があり、変更データ通知部14から伝播完了を通知されると（ステップS 1 9 3、YES）、ステップS 1 9 4としてWrite トークンを解放した後、トークン管理ノードのトークン管理部13にトークン解放成功を応答して（ステップS 1 9 5）、処理を終了する。

〔変更データ通知部〕

変更データ通知部14は、IO要求インタセプト部12または受信データ処理部15からファイルの更新データを受取り、ファイルの変更内容の他ノードへの反映をスケジュールする部分である。

【 0 1 1 6 】

変更データ通知部14は、通知されたファイルが属するオブジェクトグループの系状態テーブルに設定されている伝播モード（同期、非同期、半同期）に従い、以下の様に処理する。

【 0 1 1 7 】

尚、同期、半同期、非同期はユーザがオブジェクトグループ単位に信頼性要件に従って選択するものであるが、おおよそ以下の様な特性をもつ。

同期：ユーザプログラム17が発行したファイルへの書き込み要求の完了がユーザプログラム17に通知された時点で、ファイルへの更新データが他のノード全てに伝播されている保証が与えられる。従って、全ノードが壊れない限り、データが失われることはない。

【 0 1 1 8 】

半同期：ユーザプログラム17が発行したWrite 命令の完了がユーザプログラム17に通知された時点で、更新結果が過半数のノードに伝播している保証が与

えられる。従って、半分以上のノードが同時に壊れない限り、データが失われることはない。すなわち、ノード障害に伴う系の縮退では過半数以上のノードで新しい系を作成するので、データが失われることはない。

【0119】

非同期：ユーザプログラム17が発行したWrite 命令の完了がユーザプログラム17に通知された時点で、更新結果が他のノードに伝播している保証はない。従って、ノード障害が発生すると、完了した筈の更新結果が失われることがある。但し本実施形態のシステムでは、この場合でも更新の順序性は保証されるので、新旧のデータが入り交じって見えることはない。

1) 同期モード伝播時の処理

オブジェクトグループを構成するアクティブな全ノードに変更内容を転送し、全ノードから受信応答が戻ったところで、要求元に復帰する。

2) 半同期モード伝播時の処理

オブジェクトグループを構成するアクティブな全ノードに変更内容を転送し、過半数のノードから受信応答が全て戻ったところで要求元に復帰する。尚Write トークンは全てのノードへの伝播が完了するまでは解放しない。

3) 非同期モード伝播時の処理

変更内容をメモリ上にターゲットノード単位でキューイングし、適当なタイミングで転送する。

【0120】

ここで、適当なタイミングとは以下のいずれかの状態が発生した時を指す。

1) 系構成管理部11からSYNC要求を受け付けた時。すべての更新データを全ノードに伝播させる。

2) トークン管理部13からWrite トークンを返却する前に、FSYNC指定で呼ばれ、対象ファイルに対する変更内容を全ノードに伝播させる。

3) システムが判断した適当なタイミング。例えば一定時間立った時、あるいはキューイングされたデータが一定以上になった時。

【0121】

図23は変更データ通知部14による処理を示すフローチャートである。

変更データ通知部 1 4 は、他の構成要素から呼び出されると、まず自己を呼び出した相手を判断する（ステップ S 2 0 1）。その結果、I O 要求インタセプト部 1 2 若しくは受信データ処理部 1 5 により呼び出されたのであれば、ステップ S 2 0 2 の I O 要求インタセプト部／受信データ処理部呼び出し処理を行う。また呼び出し元が系構成管理部 1 1 であり、要求内容が S Y N C 要求であるのならば（ステップ S 2 0 3、S Y N C）、ステップ S 2 0 4 の S Y N C 要求処理を行い、また要求内容が R E S E T 要求であるのならばステップ S 2 0 5 の R E S E T 要求処理を行う。またトークン管理部 1 3 から F S Y N C 要求によって呼ばれた場合には、ステップ S 2 0 6 の F S Y N C 要求処理を行って処理を終了する。

【 0 1 2 2 】

図 2 4 は、図 2 3 のステップ S 2 0 2 の I O 要求インタセプト部／受信データ処理部呼び出し処理の動作処理を示すフローチャートである。

I O 要求インタセプト部／受信データ処理部呼び出し処理に入ると、変更データ通知部 1 4 は、呼び出し元から通知された更新要求のオブジェクトグループ番号から対応するオブジェクトグループの内部制御表を見つけ、伝播モードを調べる（ステップ S 2 1 1）。次にステップ S 2 1 2 として更新要求を更新伝播キューの最後につないでキューイングする。そして、ステップ S 2 1 1 で調べた伝播モードが非同期方式であったならば（ステップ S 2 1 3、非同期）、処理を終了し、呼び出し元に復帰する。

【 0 1 2 3 】

伝播モードが同期方式か半同期方式の場合（ステップ S 2 1 3、同期／半同期）、内部制御表の状態フラグに系再構成中が表示されていた場合には、系の再構成が完了して状態フラグの表示が消えるのを待ち合わせた後（ステップ S 2 1 4）、ステップ S 2 1 5 として系内の全アクティブノードに更新要求を送信する。

【 0 1 2 4 】

更新要求送信後、更新伝播送信キュー内の ack 待ちベクタの更新要求を送信したノードに対応するビットを立て（ステップ S 2 1 6）、応答を待つ。そして伝播モードが半同期の場合には（ステップ S 2 1 7、半同期）、ack ベクタの過半数がオフになり、更新要求を送信したノードの受信データ処理部 1 5 の過半数か

ら受信完了の応答があるまで待ち合わせ（ステップ S 2 1 8）、処理を終了して、要求元に復帰する。

【 0 1 2 5 】

また伝播モードが同期であった場合には（ステップ S 2 1 7、同期）、ステップ S 2 1 9 として ack 待ちベクタが全てオフになるのを待合わせ、トークンを自発的に返却する構成の場合トークンを解放した後に処理を終了し、要求元に復帰する。

【 0 1 2 6 】

図 2 5 は、図 2 3 のステップ S 2 0 4 の SYNC 要求処理時の変更データ通知部 1 4 の動作処理を示すフローチャートである。この SYNC 要求処理は、更新伝播送信キュー内にキューイングされている変更要求を全て系内の他ノードに伝播させて更新伝播送信キューにキューイングされている更新要求を全て掃き出させるもので、系構成管理部 1 1 から SYNC 要求により呼ばれた時に行われる。

【 0 1 2 7 】

SYNC 要求処理に入ると、変更データ通知部 1 4 は、まずステップ S 2 2 1 として内部制御表内の更新伝播送信キューのエントリを用いて更新伝播送信キューの先頭要素を読み出す。

【 0 1 2 8 】

図 2 6 は、更新伝播送信キューの構成例を示す図である。

更新伝播送信キューは、更新要求をキューイングするバッファで、内部制御表内の更新伝播送信キューエントリによって先頭要素の位置が示されるリスト構造を持つ。リスト構造の 1 つの要素は 1 つの更新要求に対応しており、変更データ通知部 1 4 は更新要求が生じると、更新伝播送信キューの最後に新規の要素を繋ぎ、処理が完了すると対応する要素を削除する。

【 0 1 2 9 】

リストデータの 1 つの要素は、次の要素の位置を示すポインタ、更新を行うファイルが属するオブジェクトグループのオブジェクトグループ番号、この更新要求を他ノードに送信したかどうかを示す送信済みフラグ、各ノード毎の応答状態を示す ack 待ちベクタ、更新対象ファイルのファイル名とそのファイル中での更

新位置をするオフセット、更新データの大きさを示す長さ、更新要求を行ったノードのノード番号を示す要求ノード番号、更新番号、依存ベクタ、更新内容を示す更新データによって構成される。これらのうち更新番号及び依存ベクタは後述する順序性保証処理で用いられるもので、順序性保証の項で詳細に説明する。

【0130】

ステップ S 2 2 1 で読み出した要素の送信済みフラグが未送信を表示していたならば (S 2 2 2、NO)、ステップ S 2 2 3 としてこの要素の更新要求を系内の全アクティブノードに送信し、送信したノードに対応するack ベクタのビットを立てる (ステップ S 2 2 4)。また読み出した要素の送信済みフラグが送信済みを表示しており、この送信要求が他ノードに伝播中のものであったならば (ステップ S 2 2 2、YES)、その要素はスキップする。

【0131】

そして次の更新伝播送信キュー内の次の要素を読みだし (ステップ S 2 2 5、YES: ステップ S 2 2 6)、ステップ S 2 2 2 ~ S 2 2 4 の処理を繰り返す。

キュー内の全要素に対して処理が完了すると (ステップ S 2 2 5、YES)、更新伝播送信キューの全要素のack 待ちベクタが0にリセットされ、更新要求を送った全てのノードから受信完了の応答があるのを待ってから (ステップ S 2 2 7)、処理を終了して要求元に復帰する。

【0132】

図 2 7 は、図 2 3 のステップ S 2 0 5 の RESET 要求処理時の変更データ通知部 1 4 の動作処理を示すフローチャートである。RESET 要求は、障害発生時に伝播途中であった要求を全ノードに伝播させ、新しい系の同期を取るなどの目的で用いられる。この RESET 要求処理は、他ノードのノード障害を認識した系構成管理部 1 1 に、RESET 要求によって呼び出された変更データ通知部 1 4 が行う処理である。RESET 要求処理では更新伝播送信キュー及び実反映遅延キューにキューイングされている更新要求を全て他ノードに伝播して更新内容を全他ノードに反映させる。

【0133】

RESET 要求処理に入ると、変更データ通知部 1 4 は、ステップ S 2 3 1 と

して図 2 6 に示した SYNC 要求処理と同様の処理を行い更新伝播送信キューにキューイングされている変更要求を全て系内の他ノードに伝播して変更内容を通知する。

【 0 1 3 4 】

次にステップ S 2 3 2 として、内部制御表内の実反映遅延キューのエントリから位置を調べ、実反映遅延キューの先頭要素を読み出す。

ステップ S 2 3 2 で読み出した要素の送信済みフラグが未送信を表示していたならば (S 2 3 3、NO)、ステップ S 2 3 4 としてこの要素の更新要求を系内の全アクティブノードに送信し、送信したノードに対応する ack ベクタのビットを立てる (ステップ S 2 3 5)。また読み出した要素の送信済みフラグが送信済みを表示しており、この送信要求が他ノードに伝播中のものであったならば (ステップ S 2 3 3、YES)、その要素はステップ S 2 3 4 及び 2 3 5 をスキップする。

【 0 1 3 5 】

そして次の実反映遅延キュー内の次の要素を読みだし (ステップ S 2 3 6、NO: ステップ S 2 3 7)、ステップ S 2 3 3 ~ S 2 3 5 の処理を繰り返す。

キュー内の全要素に対して処理が完了すると (ステップ S 2 3 6、YES)、実反映遅延キューの全要素の ack 待ちベクタが 0 にリセットされ、更新要求を送った全てのノードから受信完了の応答があるのを待ってから (ステップ S 2 3 8)、処理を終了して要求元に復帰する。

図 2 8 は、図 2 3 のステップ S 2 0 6 の F SYNC 要求処理時の変更データ通知部 1 4 の動作処理を示すフローチャートである。この F SYNC 要求処理は、変更データ通知部 1 4 が系構成管理部 1 1 からファイル名を指定して F SYNC 要求されて実行されるもので、Write トークンを解放する目的などで、更新伝播送信キュー内にキューイングされている変更要求の内指定されたファイルに対するもの全てを系内の他ノードに伝播させて更新伝播送信キューから掃き出させるものである。

【 0 1 3 6 】

F S Y N C 要求処理に入ると、変更データ通知部 1 4 は、まずステップ S 2 4 1 として内部制御表内の更新伝播送信キューのエントリを用いて更新伝播送信キューの先頭要素を読み出す。

【 0 1 3 7 】

ステップ S 2 4 1 で読み出した要素内のファイル名と指定されたファイル名と比較し同一のものであり（ステップ S 2 4 2、Y E S）、また送信済みフラグが未送信を表示していたならば（S 2 4 3、N O）、ステップ S 2 4 4 としてこの要素の更新要求を系内の全アクティブノードに送信し、送信したノードに対応するack ベクタのビットを立てる（ステップ S 2 4 5）。また要素内のファイル名が指定されたものと異なったり（ステップ S 2 4 2、N O）、ファイル名は同じであっても読み出した要素の送信済みフラグが送信済みを表示しており、この送信要求が他ノードに伝播中のものであったならば（ステップ S 2 4 3、Y E S）、その要素はスキップする。

【 0 1 3 8 】

そして次の更新伝播送信キュー内の次の要素を読みだし（ステップ S 2 4 6、Y E S：ステップ S 2 4 7）、ステップ S 2 4 2～S 2 4 5 の処理を繰り返す。

キュー内の全要素に対して処理が完了すると（ステップ S 2 4 6、Y E S）、更新伝播送信キューのステップ S 2 4 5 でビットを立てたack 待ちベクタが 0 にリセットされ、更新要求を送った全てのノードから受信完了の応答があるのを待ってから（ステップ S 2 4 8）、処理を終了して要求元に復帰する。

【 0 1 3 9 】

その後変更データ通知部 1 4 は、適当なタイミングで実反映遅延キューを先頭からスキャンし、まだ他ノードに伝播されていない先頭から特定数の変更要求を全アクティブノードに転送する。

〔受信データ処理部〕

受信データ処理部 1 5 は、他ノードからデータを受信し、自ノードへの反映処理を行う部分である。

【 0 1 4 0 】

受信データ処理部 1 5 が他ノードから受取るデータには、Read/Write要求、R

RESET要求及び等価性回復転送データの4種類があり、受信データ処理部15はそれぞれに応じた処理を行う。

【0141】

図29は、受信データ処理部15の動作処理を示すフローチャートである。

受信データ処理部15は、他ノードから要求を受信すると、まずその内容を判断する(ステップS251)。その結果更新要求であれば、ステップS252の更新要求処理を行う。また自ノードがWrite トークンを保持しており、他ノードからRead要求若しくはWrite 要求が送信されてきたのであれば、ステップS253のRead/Write要求処理を行う。また、他ノードが離脱したノードを検出してRESET要求を送信してきたのならば、ステップS254のRESET要求処理を行う。また、自ノードが等価性回復中で、等価性回復転送要求をしたノードから等価性回復転送データを送信してきたのならば、ステップS255の等価性回復転送データ処理を行う。

【0142】

図30は、図29のステップS252の更新要求処理における受信データ処理部15の処理を示すフローチャートである。

更新要求処理に入ると、受信データ処理部15は、受信した更新データに対応するオブジェクトグループの内部制御表を参照し、このオブジェクトグループの伝播モードと等価性回復中であるかどうかを調べる。その結果伝播モードが同期モードあるいは半同期モードであるか(ステップS261、YES)、非同期モードであっても状態フラグに等価性回復中が表示されていた場合(ステップS261、NO:ステップS262、YES)、OSファイルシステムを介して、自ノードの対応ファイルに変更データを直ちに反映させ(ステップS263)、処理を終了する。

【0143】

また転送モードが非同期モードであり(ステップS261、NO)、また等価性回復中でなかった場合には(ステップS262、NO)、ステップS262として受信した変更要求を実反映遅延キューの最後尾に繋ぎ、順序性保証を考慮して変更要求を自ファイルへ反映させる。尚順序性保証については後述する。

【 0 1 4 4 】

図 3 1 は、実反映遅延キューの構成例を示す図である。

実反映遅延キューは、非同期モードによる更新要求をキューイングするバッファで、内部制御表内の実反映遅延キューエントリによって先頭要素の位置が示されるリスト構造を持つキュー部分 2 1 と受信済みベクタ 2 2 によって構成される。キュー部分 2 1 の 1 つの要素は 1 つの更新要求に対応しており、受信データ処理部 1 5 は、非同期モードのオブジェクトグループ内のファイルに対する更新要求を受信すると、実反映遅延キューの最後に新規の要素を繋ぎ、処理が完了すると対応する要素を削除する。

【 0 1 4 5 】

キュー部分 2 1 の 1 つの要素は、基本的に更新伝播送信キューの要素と同じ構成で、次の要素の位置を示すポインタ、更新を行うファイルが属するオブジェクトグループのオブジェクトグループ番号、この更新要求を他ノードに送信したかどうかを示す送信済みフラグ、各ノード毎の応答状態を示すack 待ちベクタ、更新対象ファイルのファイル名とそのファイル中での更新位置をするオフセット、更新データの大きさを示す長さ、更新要求を行ったノードのノード番号を示す要求ノード番号、更新番号、依存ベクタ、更新内容を示す更新データによって構成される。

【 0 1 4 6 】

これらのうち更新番号及び依存ベクタは後述する順序性保証処理で用いられるもので、順序性保証の項で詳細に説明する。また送信済みフラグ及びack 待ちベクタは、系構成管理部 1 1 から R E S E T 要求を受けた時にのみ用いられる。

【 0 1 4 7 】

また受信済みベクタ 2 2 は、系内のノード分の要素を備え受信した更新要求内の依存ベクタ最新の依存ベクタが記録される。尚この点についても、順序性保証の項で詳細に説明する。また受信済みマトリックスについても順序性保証の項で説明する。

【 0 1 4 8 】

図 3 2 は、図 2 9 のステップ S 2 5 3 のRead/Write要求処理における受信デー

タ処理部 1 5 の処理を示すフローチャートである。

Read/Write要求処理に入ると受信データ処理部 1 5 の処理は、受信したRead要求若しくはWrite 要求にオプションで F O R C E が指定されているかどうかによって処理が異なる。

【 0 1 4 9 】

受信したRead/Write要求が、等価性回復中のノードからのものであり、 F O R C E オプションの指定されたものである時は（ステップ S 2 7 1、 Y E S）、ステップ S 2 7 2 としてトークン管理部 1 3 に要求処理に必要なReadトークン若しくはWrite トークンの獲得を依頼する。その結果獲得に成功すれば（ステップ S 2 7 3、 Y E S）、ステップ S 2 7 4 に処理を移し、獲得に失敗すれば（ステップ S 2 7 3、 N O）、ステップ S 2 7 8 として要求元ノードにエラー応答を行った後処理を終了する。

【 0 1 5 0 】

また受信したRead/Write要求が、 F O R C E オプションの指定の無いものである時は（ステップ S 2 7 1、 N O）、自ノードがWrite トークンを保持していない時は（ステップ S 2 7 9、 N O）、ステップ S 2 7 8 として要求元ノードにエラー応答を行った後処理を終了する。また自ノードがWrite トークンを保持している時は（ステップ S 2 7 9、 Y E S）、ステップ S 2 7 4 に処理を移す。

【 0 1 5 1 】

ステップ S 2 7 4 では、内部制御表を参照し、Read/Write要求の対象となっているオブジェクトグループの伝播モードを調べる。その結果同期モードあるいは半同期モードであった場合には（ステップ S 2 7 4、同期／半同期）、O S のファイルシステムに依頼して、要求された処理を行い（ステップ S 2 7 6）、結果を要求もとノードに応答して処理を終了する。尚ステップ S 2 7 6 において、Write 要求に対する処理の場合、自ファイルへの書込み処理の他、変更内容の他ノードへの伝播を変更データ通知部 1 4 に依頼する。

【 0 1 5 2 】

ステップ S 2 7 4 で、Read/Write要求の対象となっているオブジェクトグループの伝播モードが非同期であるならば（ステップ S 2 7 4、非同期）、後述する

順序性保証の項で述べる順序性保証の為の処理を考慮しつつ、I O 要求インタセプト部 1 2 による Read/Write 要求処理に準じた処理を行い、結果を要求元ノードに返し（ステップ S 2 7 7）、処理を終了する。

図 3 3 は、図 2 9 のステップ S 2 5 4 の R E S E T 要求処理における受信データ処理部 1 5 の処理を示すフローチャートである。

【 0 1 5 3 】

R E S E T 要求処理に入ると、受信データ処理部 1 5 は、ステップ S 2 8 1 として内部制御表内の実反映遅延キューのエントリから位置を調べ、実反映遅延キューの先頭要素を読み出す。そしてその要素が、系から離脱したノードからの更新要求を待っているものであるならば（ステップ S 2 8 2、Y E S）、ステップ S 2 8 3 としてその更新要求を実反映遅延キューから削除して解放する。また他のノードからの更新要求であったならばそのまま残しておく（ステップ S 2 8 2、N O）。

【 0 1 5 4 】

そして次の実反映遅延キュー内の次の要素を読みだし（ステップ S 2 8 4、N O：ステップ S 2 8 5）、ステップ S 2 8 2 ～ 2 8 4 の処理を繰り返し、キュー内の全要素に対して処理が完了すると（ステップ S 2 8 4、Y E S）、処理を終了する。

【 0 1 5 5 】

図 3 4 は、図 2 9 のステップ S 2 5 5 の等価性回復データ処理における受信データ処理部 1 5 の処理を示すフローチャートである。

等価性回復データ処理に入ると、ステップ S 2 9 1 として受信データ処理部 1 5 は、ステップ S 2 9 1 としてファイルシステムを呼び出し、受信した等価性回復転送データの自ノードのファイルへの反映を依頼し、完了応答を待った後（ステップ S 2 9 2）、処理を終了する。

〔順序性保証〕

本システムでは、ファイルの更新を行うと更新内容は更新要求として系内の他ノードに伝播されてゆく。伝播モードとしては、同期、非同期、半同期の 3 つの

モードがあり、このうち同期モード及び半同期モードによる伝播以外の時は、系の縮退時に完了した筈のファイルの更新の結果が失われてしまう可能性がある。この為、系縮退時に一部データが失われ、結果として新旧データが入り乱れる事態が生じる。半同期モードではしかもファイルへの更新データが他のノードに更新された順番に届くとは限らない。

【 0 1 5 6 】

本実施形態では、非同期モード時、受信した更新データを実反映遅延キューにキューイングしてゆき、実反映遅延キュー内の更新データの自ファイルへの反映を更新番号と依存ベクタによって管理し、順序性保証を行い、系縮退時に新旧データが入り乱れることを防止する。

【 0 1 5 7 】

この更新番号と依存ベクタは、例えば内部制御表内に設定される。内部制御表は、オブジェクトグループ毎に展開されるので、この構成の場合、高新番号と依存ベクタもオブジェクトグループ毎に持つことになる。従ってオブジェクトグループを互いに関係があるファイルのみで定義すれば、互いに無関係な更新間の順序性保証は行われず、オーバーヘッドを削減することが出来る。

1) 更新番号

更新番号は、系内で発生するファイル更新のノード内に閉じた順序性を表す為に単調に増加する番号でありノード毎、オブジェクトグループ毎に用意する。I/O要求インタセプト部12はユーザプログラムからWrite要求を受ける度にこの更新番号をインクリメントして更新する。

2) 依存ベクタ

依存ベクタは、他ノードの更新番号を含むベクタで、「更新番号で示される更新要求が依存する」他ノードが行った更新を特定する。依存ベクタは、オブジェクトグループ毎に用意され、そのオブジェクトグループに属するノード数個の要素をもつ。

【 0 1 5 8 】

各要素の内、自ノードに対応する部分には、常に自ノードの更新番号より1つ小さい値が設定される。依存ベクタは、更新データの伝播時に、更新番号共に更

新データに付加されて伝播される。

【 0 1 5 9 】

Write トークンの獲得に失敗してWrite 処理を他ノードに依頼する場合、I O 要求インタセプト部 1 2 がWrite 要求に更新番号と依存ベクタを付加し、これを依頼先のノードに送信する。Write 要求によるファイルの更新内容は、Write 要求を受けたノード経由で更新伝播時に系内の全ノードに通知される。

【 0 1 6 0 】

またRead要求を依頼されたノードは、応答も依存ベクタを付加する。

図 3 5 は、Write 要求及びRead要求の応答に付加される依存ベクタの例を示す図である。

【 0 1 6 1 】

同図上段は 3 つのノードで系が構成されている場合に、Write 要求ノード 2 からノード 1 にWrite 要求を行う場合を示す図であり、下段はRead要求に対する応答を行う場合を示している。

【 0 1 6 2 】

ノード 2 の I O 要求インタセプト部 1 2 は、ユーザプログラム 1 7 からWrite 要求を受けると、内部制御表内の更新番号及び依存ベクタ内の自己に対応する部分をインクリメントし（同図の場合更新番号を 9 → 1 0、依存ベクタを（1 0、8、6）→（1 0、9、6）に変更）、これを更新番号と共にWrite 要求に付加してノード 1 に送る。またRead要求の応答の場合にはこれらの更新は行わず、内部制御表に設定されている依存ベクタをそのまま付加して送信する。

【 0 1 6 3 】

ノード 1 では、Write 要求の場合受信した更新データを更新番号及び依存ベクタと共に実反映遅延キューにキューイングすると共に、内部制御表内の依存ベクタと各要素毎に受信済みベクタ 2 2 内のベクタと受信した依存ベクタとを比較（ノード 2 の部分は更新番号と比較）、受信したベクタの方が大きければこれを新たな値として内部制御表にセットする。

【 0 1 6 4 】

また図 3 5 下段は、Read要求の応答に対しては、単に内部制御表の内の依存ベ

クタと応答に付加していた依存ベクタとを要素毎に比較し、受信したベクタの方が大きければこれを新たな値として内部制御表にセットする。

【0165】

依存ベクタは、他ノードから送信されてきた更新要求や、Write 要求で通知された更新データを実ファイルに反映してもよいかどうかを受信データ処理部15が判断するのに使用する。受信データ処理部15は、依存ベクタ内の要素全てのより小さい更新番号の更新要求を各ノードから全て受取済の場合には、実ファイルに反映してよいと判断して更新を行う。

【0166】

尚受信した更新要求より先行する更新に対する更新要求にまだ到着していないものが存在する場合、系再構成時の破棄に備え、その未着の更新内容が送られて来るまで受信した更新内容を実反映遅延キューに保持しておき、実ファイルへの反映を遅らせる。これにより、更新内容が前後して届いた場合に系の再構成が生じて、データが破壊されることはない。

【0167】

図36は、受信データ処理部15が行う依存ベクタによる判断処理を説明する図である。

同図はノード3の実反映遅延キューの状態を示したもので、キューには受信順にノード1からの更新番号12の更新要求（同図中1/12）、ノード1からの更新番号13の更新要求（同図中1/13）及びノード2からの更新番号12の更新要求（同図中2/12）更新要求が実反映遅延キューにキューイングされている。また受信済みベクタ22から、既に反映済みの更新データとして更新番号がノード1及び2は更新番号10まで、ノード3は更新番号5までの更新データが自ファイルに反映されていることが判る。

【0168】

この状態を初期状態T0とし、次の状態T1としてノード2から更新番号の11の更新要求（依存ベクタ（10, 10, 5））がノード3に到着したとする。

これにより、ノード2からの更新要求は更新番号がまで全て12揃ったことになるので（受信済みベクタ22から更新番号10以前のものは既に反映済み）、

受信済みベクタ 2 2 を (1 0 , 1 0 , 5) から (1 0 , 1 2 , 5) と変更すると共に反映可能となった 2 / 1 1 の更新データを自ファイルに反映させる。しかし、2 / 1 2 の更新データに関しては、2 / 1 2 の更新データの依存ベクタと受信済みベクタ 2 2 の値とを比較すると、ノード 1 の部分の値が 2 / 1 2 の更新要求の方が大きいので、これは自ファイルには反映させずに実反映遅延キュー内に保持しておく。

【 0 1 6 9 】

また次の T 3 の状態として、ノード 1 から更新番号 1 1 の要求 (要求 1 / 1 1 (1 0 , 1 1 , 5) が到着したとする。これによりノード 1 からの更新要求は更新番号 1 3 まで全てノード 3 に到着したことになるので、受信済みベクタ 2 2 を (1 0 , 1 2 , 5) から (1 3 , 1 2 , 5) に変更すると共に、反映可能となった要求 1 / 1 1 , 1 / 1 2 , 1 / 1 3 , 2 / 1 1 を全て実ファイルに反映させ、これらを実反映遅延キューから削除する。

【 0 1 7 0 】

また Read 要求を処理する場合では、実反映遅延キュー対応するデータが退避されていればそちらを優先して読みだし、要求元に送る。この際、応答する依存ベクタもキューイングされているデータに付加されているものを返す。

【 0 1 7 1 】

この様に処理することにより更新要求が実際の更新順から前後して届いても、受信データ処理部 1 5 は、順序性保ったデータの更新を行うことが出来る。

尚、実反映遅延キューからデータを返す処理を不要にして制御を単純化するために、Write 要求を受取った受信データ処理部 1 5 が、Write 要求に付加された依存ベクタからその Write 要求に依存関係の有るデータが自ノードに全て到着するのを待合わせる構成としても良い。この場合 Write 要求された更新データの自ファイルへの反映と Write トークンの解放をその Write トークン下で行い、更新を依存するデータが全ノードに到着したことを確認出来るまで遅らせる。この点については後述する。

【 0 1 7 2 】

この構成の場合、自ノードのデータを Read 使用とする場合には、Write トーク

ンの解放を介して、依存するデータが自ノードに反映済みとなるので、Read要求の処理で実反映遅延キューからデータを取り出し応答するという処理が不要となる。ただしこの場合でも、系再編成によりデータの順序性が崩れることを防ぐため、実反映遅延キューを介して、実ファイルのへ反映を遅らせる処理は依然必要となる。

3) 依存ベクタの更新タイミング

依存ベクタは以下のタイミングで更新される。

a) 他ノードからWrite 要求が送られてきた時

受信データ処理部 1 5 は自ノードの依存ベクタの要求元ノードに対応する要素に送られてきた更新番号を設定する。

b) I O要求インタセプト部 1 2 が他ノードにRead要求を送り、応答としてReadデータをもらった時

受信データ処理部 1 5 は、応答と共に送られてきた依存ベクタと自身が内部制御表内に保持する依存ベクタとを要素毎に比較し、大きい値を内部制御表内に設定する。Read要求を受けたときは受信データ処理部 1 5 は、Read要求を受けた時点の依存ベクタを応答に付加して返す。

【 0 1 7 3 】

上記の様に依存ベクタを伝播することで、複数のノード間に跨がるデータ間の依存性を表現することが出来る。例えば、a (ノード 1) → b (ノード 2) → c (ノード 3) で表現される依存関係がある更新要求の場合、ノード 3 から送られてきた更新要求 c は更新要求 a、b の更新が伝播するまで不揮発化が延ばされる。

【 0 1 7 4 】

図 3 7 は、依存関係のある更新要求の順序性の保証を示す図である。

同図は同一のオブジェクトグループに属するファイル f a、f b 及び f c の 3 のファイルに対し 3 つのノード 1、2 及び 3 によってRead/Write要求が発生した場合の依存ベクタによるを示したもので、t 0 ~ t 5 の順でファイルに対する更新が行われた場合、t 0、t 2、t 4 の 3 つの状態が発生した更新要求に付加される依存ベクタには、 $(0, 0, 0) < (1, 0, 0) < (1, 1, 0)$ の関係が有るので、各ノードに更新要求が順不同で届いてもファイルには順番に反映さ

れる。

4) 参照要求時

ユーザプログラム 17 からの Read 要求に対し、他のノードに Read 要求を依頼して応答結果を得る場合、I/O 要求インタセプト部 12 は、受取ったデータに付加されている依存ベクタで示された受信データに依存関係が有る更新要求を全て受信するまで、ユーザプログラム 17 に参照結果を渡さない。

【0175】

この様にユーザプログラム 17 に応答を返すのを遅らせて、同期を取ることで、系の再構成を跨がってこのノードが生き続けた場合に、ユーザプログラム 17 が参照したデータが失われてユーザプログラム 17 の誤動作を防ぐことが出来る。

【0176】

尚処理を単純にするため、他ノードに Read 要求に対する応答を返す場合、受信データ処理部 15 で、自ノードがそれまでに行った変更が過半数のノードに伝わるのを待ってから応答を返すという構成にすることも出来る。この構成の場合には、他ノードに Read 要求の応答結果を返す時にはその応答結果が依存する更新要求が系内の過半数のノードに必ず反映済みであることが保証される。よって、更新要求 a (ノード 1) → 更新要求 b (ノード 2) → 更新要求 c (ノード 3) の様な間接的な依存関係が更新に対しても、ノード 2 がノード 1 から Read データを受信した時点で、更新要求 a が過半数のノードに伝播していることになり、ノード 3 がノード 2 から Read 結果を受信した時点では依存関係にある更新要求 a が過半数に伝播していることが保証されることになる。

【0177】

更に、図 31 に示す受信済みマトリックスを導入して、WRITE トークンの回収を WRITE トークンで保護された更新と依存関係の有る更新要求が全ノードに伝わるまで遅らせる最適化を行う構成とすることも可能である。

【0178】

この構成の場合、更新要求は依存関係を持つ更新要求が系内の全ノードに伝播されるまで更新伝播送信キューに繋がれたままとなる。抛って、Read 要求に対し

更新伝播送信キューに繋がっていないデータを返す場合には、依存するデータが既に系内の全ノードに伝わっている保証がとれる。

【0179】

従って、他ノードからのRead要求に対し、要求を依頼されたノードは更新伝播送信キューにあるデータを応答とするときのみ、その応答としたデータに対応する依存ベクタを応答すればよく、更新伝播送信キューにないデータを応答とする場合には、依存ベクタなしを応答することができる。依存ベクタなしを応答されたREAD要求ノードは依存関係に変更がないので自身の依存ベクタを更新したり、依存ベクタで規定される更新要求が到着するのを待ち合わせる必要がなくなる。

【0180】

図31に示す受信済みマトリックスは、ノード毎に存在するマトリックスで、他ノードの受信済みベクタを要素として持ち、自ノードが認識している他ノードの進行状況を示す。上記したWRITEトークンの回収をWRITEトークンで保護された更新と依存関係の有る更新要求が全ノードに伝わるまで遅らせる構成の場合、WRITEトークン保持ノードは、この受信済みマトリックスから、依存関係の有る更新が全ノードに伝達されたことを認識する。

【0181】

各ノードは一定時間毎に、系内の全ての他ノードに対し自身の受信済みマトリックスをメッセージとして広報し、このメッセージを受信したノードは自身の受信済みマトリックスを更新する。受信マトリックスの更新方法是对应する受信済みベクタに対し、依存ベクタの更新方法で説明したのと同じ方法を適用すればよい。

5) データ更新時

他ノードにWrite 要求を依頼した場合、IO要求インタセプト部12は応答で通知される依存ベクタ（更新伝播送信キューに存在する同一ファイルに対する更新の最終要求を示す依存ベクタ）からそれ以前の更新における更新データが全て到着するのを待合わせ、その後自身のデータも更新する。

【0182】

Write 要求は自身がそれ以前に行ったRead/Write要求に依存している。このうち、自身のWrite データは上記待合わせ処理により自身で反映済みであることが保証される。

【0183】

また、参照データに関しては4)で述べた処理により、受取ったデータが依存するデータが全て自ノードに反映済みであることが保証される。従って、Write 要求時点で更新要求のデータが依存する他の更新データが自ノードでのファイルに反映済みである保証が得られている。尚更新データを他ノードからの伝播を待たず自ノードに反映しておくのは4)で述べたのと同じ理由で系再編を跨がって動作を続けるユーザプログラム17の誤動作を防止するためである。

【0184】

一方更新データの自ノードへの反映を先に行うと、同じファイルに対する古い伝播が後で到着したり、その更新が前提とする更新が系再編で失われることがある。この事態を防ぐために応答で通知された依存ベクタの最大のものを使い、依存関係のある更新データを待合わせる必要がある。

【0185】

図38は他ノードのWrite 要求を処理する時において、更新伝播送信キューに同じファイルに対する更新要求が存在していた場合の処理を説明する図である。

更新伝播送信キューが同図の状態、ファイルf aに対するWrite 要求を受けると、受信データ処理部15は同じファイルf aに対する最遅の更新要求（要求2/12）に対応する（11, 12, 6）を依存ベクタとして応答する。もし、更新伝播遅延キューに同じファイルに対する要求が存在しなければ、依存ベクタ無しを応答する。

【0186】

図39は、本実施形態における上記ファイルレプリケーション制御をコンピュータプログラムにより実現した場合の各ノードの構成を示す図である。

各ノードは図39の様にCPU31、ROM、RAMによる主記憶装置32、補助記憶装置33（図4のローカルディスク装置に対応）、ディスプレイ、キーボード等の入出力装置（I/O）34、LANやWAN、一般回線等により他ノ

ードとネットワーク接続を行うモデム等のネットワーク接続装置 3 5 及びディスク、磁気テープなどの可搬記録媒体 3 7 から記憶内容を読み出す媒体読取り装置 3 6 を有し、これらが互いにバス 3 8 により接続される構成を備えている。

【 0 1 8 7 】

また図 3 9 の情報処理システムでは、媒体読取り装置 3 6 により磁気テープ、フロッピーディスク、CD-ROM、MO 等の記録媒体 3 7 に記憶されているプログラム、データを読み出し、これを主記憶装置 3 2 またはハードディスク 3 3 にダウンロードする。そして本実施形態による各処理は、CPU 3 1 がこのプログラムやデータを実行することにより、ソフトウェア的に実現することが可能である。

【 0 1 8 8 】

また、このノードでは、フロッピーディスク等の記録媒体 3 7 を用いてアプリケーションソフトの交換が行われる場合がある。よって、本発明は、計測制御用監視端末やその画面の操作方法に限らず、コンピュータにより使用されたときに、上述の本発明の実施の形態の機能をコンピュータに行わせるためのコンピュータ読み出し可能な記録媒体 3 7 として構成することもできる。

【 0 1 8 9 】

この場合、「記録媒体」には、例えば図 4 0 に示されるように、CD-ROM、フロッピーディスク（あるいは MO、DVD、リムーバブルハードディスク等であってもよい）等の媒体駆動装置 4 7 に脱着可能な可搬記録媒体 4 6 や、ネットワーク回線 4 3 経由で送信される外部の装置（サーバ等）内の記憶手段（データベース等）4 2、あるいは情報処理装置 4 1 の本体 4 4 内のメモリ（RAM 又はハードディスク等）4 4 等が含まれる。可搬記録媒体 4 6 や記憶手段（データベース等）4 2 に記憶されているプログラムは、本体 4 4 内のメモリ（RAM 又はハードディスク等）4 5 にロードされて、実行される。

【 0 1 9 0 】

【発明の効果】

本発明によれば、共用ファイルへのアクセス要求が生じたノードに対し、その共用ファイルに対する最新のデータを保持するノードが通知される。よって、共

用ファイルをアクセスするノードは常に最新のデータに対してアクセスすることが出来る。また各ノードは同一のデータを参照することになるので、各ノードからは一貫性の有るデータが見える。

【 0 1 9 1 】

また各ノードは、トークンの獲得に失敗してもトークンを獲得できるまで待つことなく処理を続行できる。更に複数のノードによる同一のファイルに対する同時アクセスを可能としている。この為、高い反応性を持つシステムを構築することができる。

ータが見える。

【 0 1 9 2 】

更に更新内容を他ノードに非同期で伝送しても、全ノードから同じデータが見える。

また更新データには更新の順序性、依存性を示す情報が付加されており、この情報に基づいてファイルの更新が行われるので、途中で系の再構成が生じても、データ更新の順序性が壊れることはない。また動作中の他ノードから矛盾したデータが見えることはない。

【 0 1 9 3 】

更に、1乃至複数のファイル毎に更新内容の伝播方式や伝播させるノードを設定できるので、業務の性格や性能要件に基づいて設定を行える。

また、新規ノードの参加時において、最新データの復元処理中に生じたアクセス要求を最新データを保持している他のノードに送ることにより、復元処理の完了を待たずに新規参加ノードの業務を開始することが出来る。更に、この時、系内で復元処理と平衡して現在系に加わっているノードの業務を続行できる。

【 0 1 9 4 】

又共用ファイルに対する処理を、該共用ファイルを共用する他ノードと同期して停止する整然停止を行った場合、共用ファイルへの処理を再開する際、他ノードと同期して再開することにより共用ファイルに対するデータの復元処理を行う必要が無い。

【図面の簡単な説明】

【図 1】

本発明の原理図である。

【図 2】

系の構成を示す図である。

【図 3】

本発明における基本原理を示す図である。

【図 4】

本実施形態の系を構成するノードの構成を示すブロック図である。

【図 5】

系状態テーブルの構成例を示す図である。

【図 6】

内部制御表の構成例を示す図である。

【図 7】

Joinコマンド投入時の系構成管理部による動作処理を示すフローチャートである。

【図 8】

参入処理時の系構成管理部の動作処理を示すフローチャートである。

【図 9】

JOIN要求受付処理時の系構成管理部の動作処理を示すフローチャートである。

【図 1 0】

Join通知を受取ったノードの系構成管理部が行う処理を示すフローチャートである。

【図 1 1】

等価性回復処理の系構成管理部の動作処理を示すフローチャートである。

【図 1 2】

等価性回復転送要求を受信したノードの系構成管理部が行う動作処理を示すフローチャートである。

【図 1 3】

等価性回復完了メッセージを受信したノードの系構成管理部が行う動作処理を示すフローチャートである。

【図 1 4】

leave コマンドを投入された時の時の系構成管理部の動作処理を示すフローチャートである。

【図 1 5】

系内の他ノードの離脱を認識したノードの系構成管理部の処理動作を示すフローチャートである。

【図 1 6】

I O 要求インタセプト部による処理動作を示すフローチャートである。

【図 1 7】

トークン制御表の構成例を示す図である。

【図 1 8】

トークン管理ノードのトークン管理部の処理動作を示すフローチャートである。

【図 1 9】

Write トークン獲得要求処理時のトークン管理部の処理動作を示すフローチャートである。

【図 2 0】

Read トークン獲得要求処理時のトークン管理部の処理動作を示すフローチャートである。

【図 2 1】

トークン解放／回収要求処理時のトークン管理部の処理動作を示すフローチャートである。

【図 2 2】

不要になったトークンを自発的に返却しない構成の場合に発行されるWrite トークン回収要求を受けたWrite トークン保持ノードが行う動作処理を示すフローチャートである。

【図 2 3】

変更データ通知部による動作処理を示すフローチャートである。

【図 2 4】

I O 要求インタセプト部／受信データ処理部呼び出し処理の変更データ通知部の動作処理を示すフローチャートである。

【図 2 5】

S Y N C 要求処理時の変更データ通知部の動作処理を示すフローチャートである。

【図 2 6】

更新伝播送信キューの構成例を示す図である。

【図 2 7】

R E S E T 要求処理時の変更データ通知部の動作処理を示すフローチャートである。

【図 2 8】

F S Y N C 要求処理時の変更データ通知部の動作処理を示すフローチャートである。

【図 2 9】

受信データ処理部の動作処理を示すフローチャートである。

【図 3 0】

更新要求処理における受信データ処理部の動作処理を示すフローチャートである。

【図 3 1】

実反映遅延キューの構成例を示す図である。

【図 3 2】

Read/Write 要求処理における受信データ処理部の処理を示すフローチャートである。

【図 3 3】

R E S E T 要求処理における受信データ処理部の動作処理を示すフローチャートである。

【図 3 4】

等価性回復データ処理における受信データ処理部 1 5 の動作処理を示すフローチャートである。

【図 3 5】

Write 要求及びRead要求の応答に付加される依存ベクタの例を示す図である。

【図 3 6】

依存関係のある更新要求の順序性の保証を示す図である。

【図 3 7】

依存関係のある更新要求の順序性の保証を示す図である。

【図 3 8】

Write 要求を自ノードで処理する時において、実反映遅延キューに同じファイルに対する更新要求が存在していた場合の処理を説明する図である。

【図 3 9】

ノードとなる計算機システムの環境図である。

【図 4 0】

記憶媒体の例を示す図である。

【符号の説明】

A ～ J ノード

- 1 1 系構成管理部
- 1 2 I O 要求インタセプト部
- 1 3 トークン管理部
- 1 4 変更データ通知部
- 1 5 受信データ処理部
- 2 1 キュー部分
- 2 2 受信済みベクタ
- 3 1 C P U
- 3 2 主記憶装置
- 3 3 補助記憶装置
- 3 4 入出力装置
- 3 5 ネットワーク接続装置

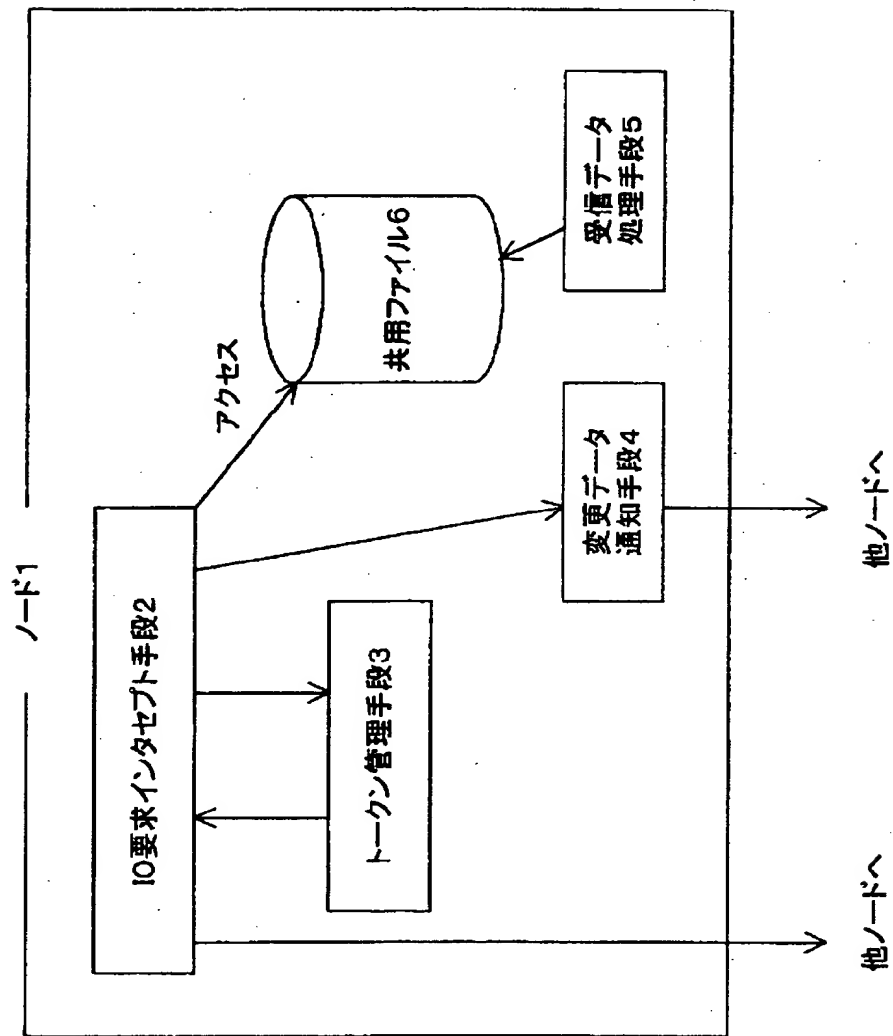
- 3 6 媒体読取り装置
- 3 7 可搬記憶媒体
- 3 8 バス
- 4 1 情報処理装置
- 4 2 記憶手段
- 4 3 ネットワーク回線
- 4 4 情報処理装置本体（コンピュータ）
- 4 5 メモリ
- 4 6 可搬記録媒体

【書類名】

図面

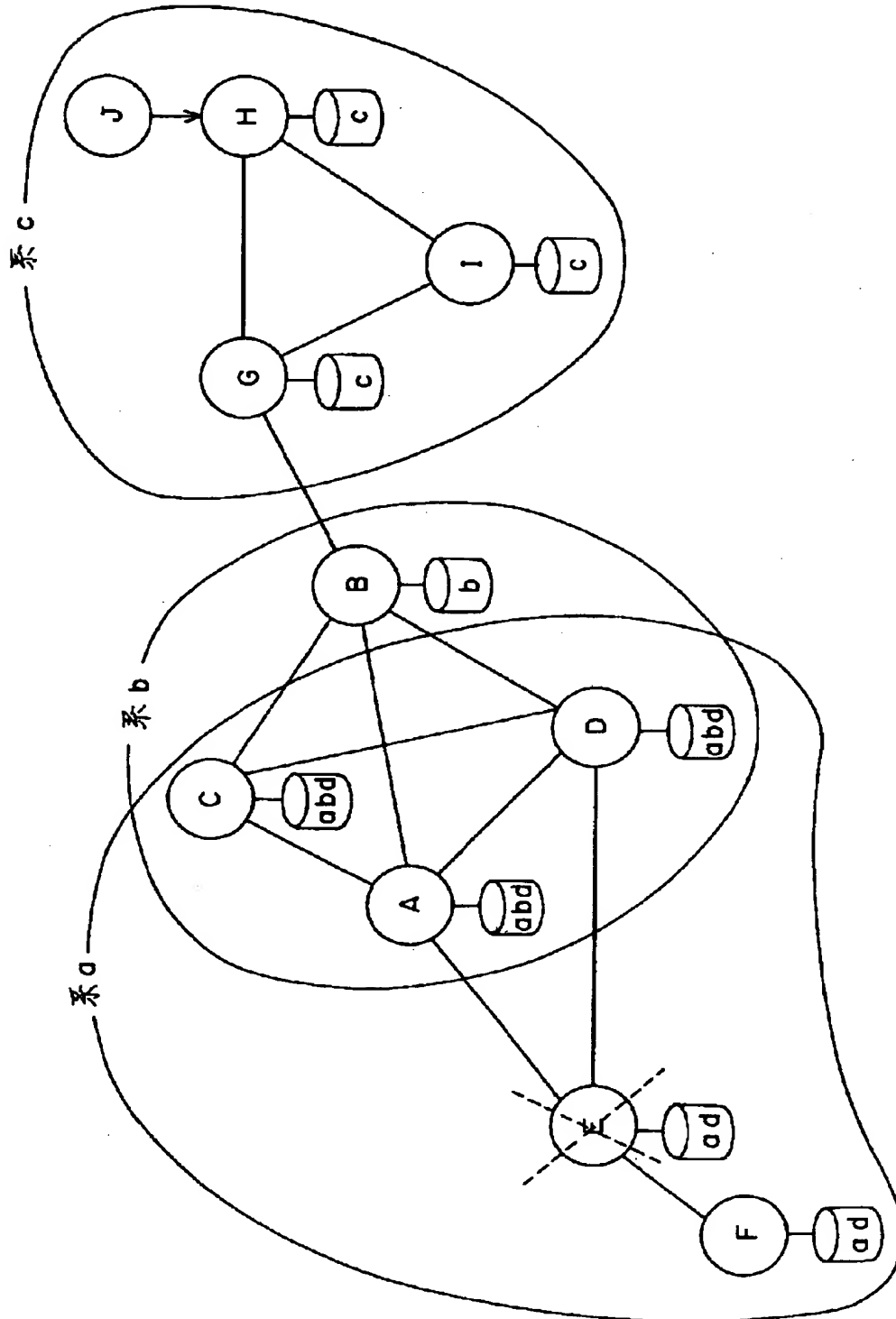
【図 1】

本発明の原理図



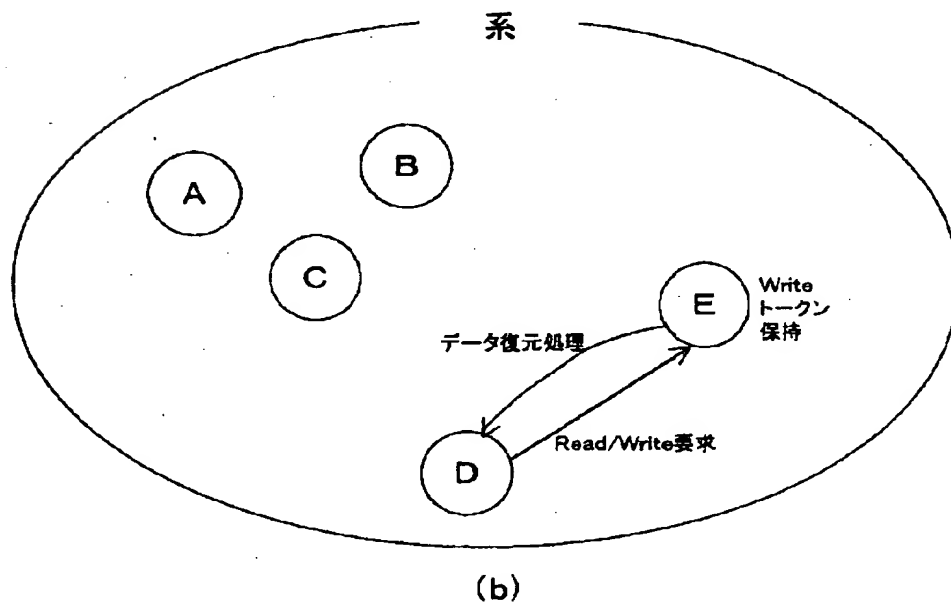
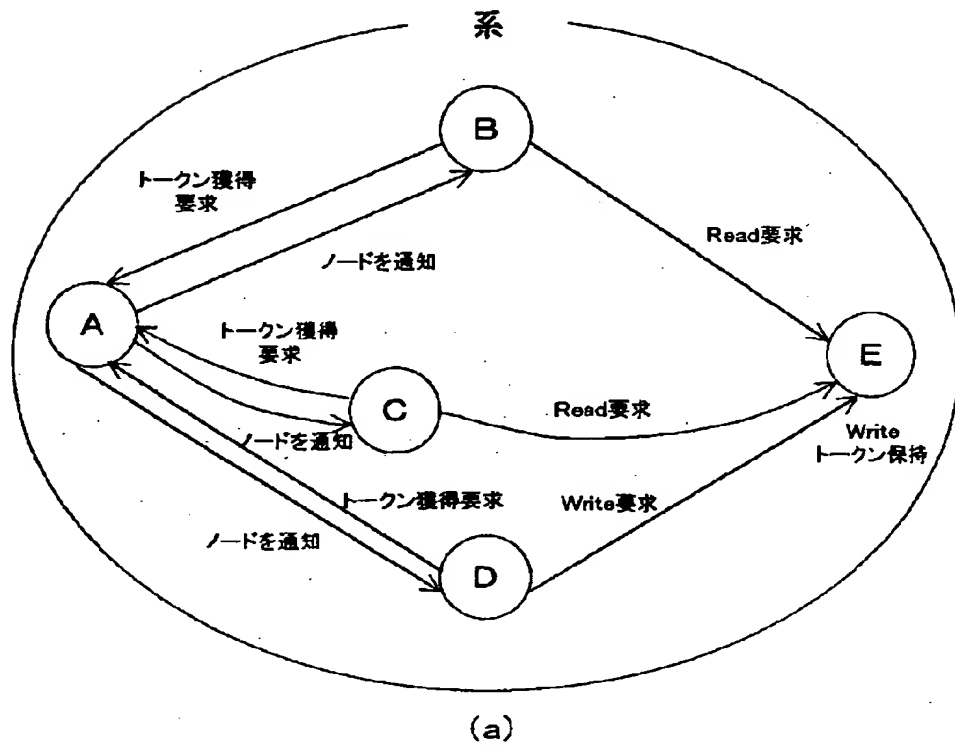
【図2】

系の構成を示す図



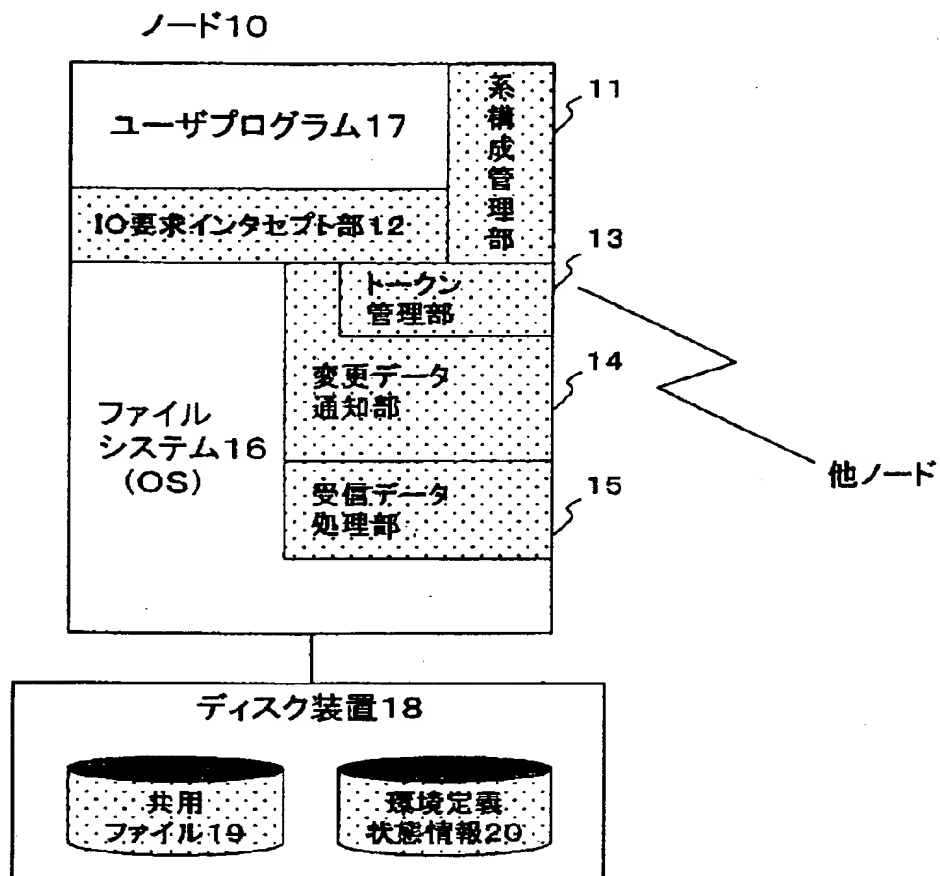
【図3】

本発明における基本原理を示す図



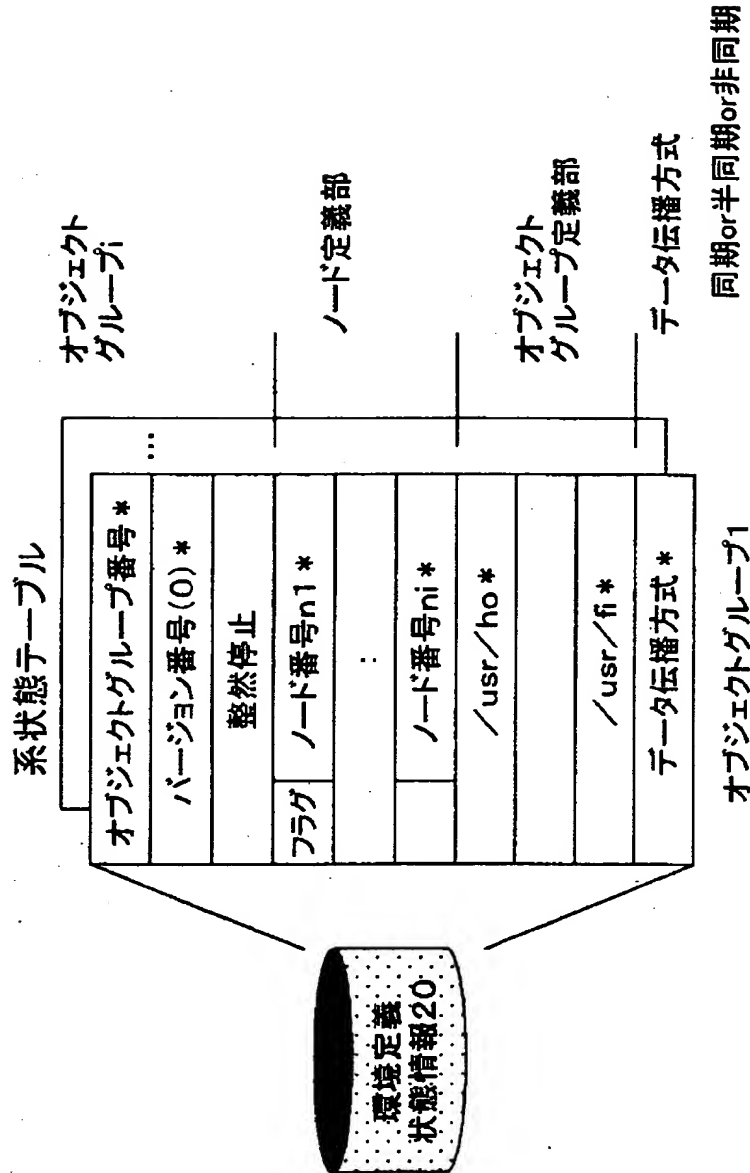
【図 4】

本実施形態の系を構成するノードの構成を示すブロック図



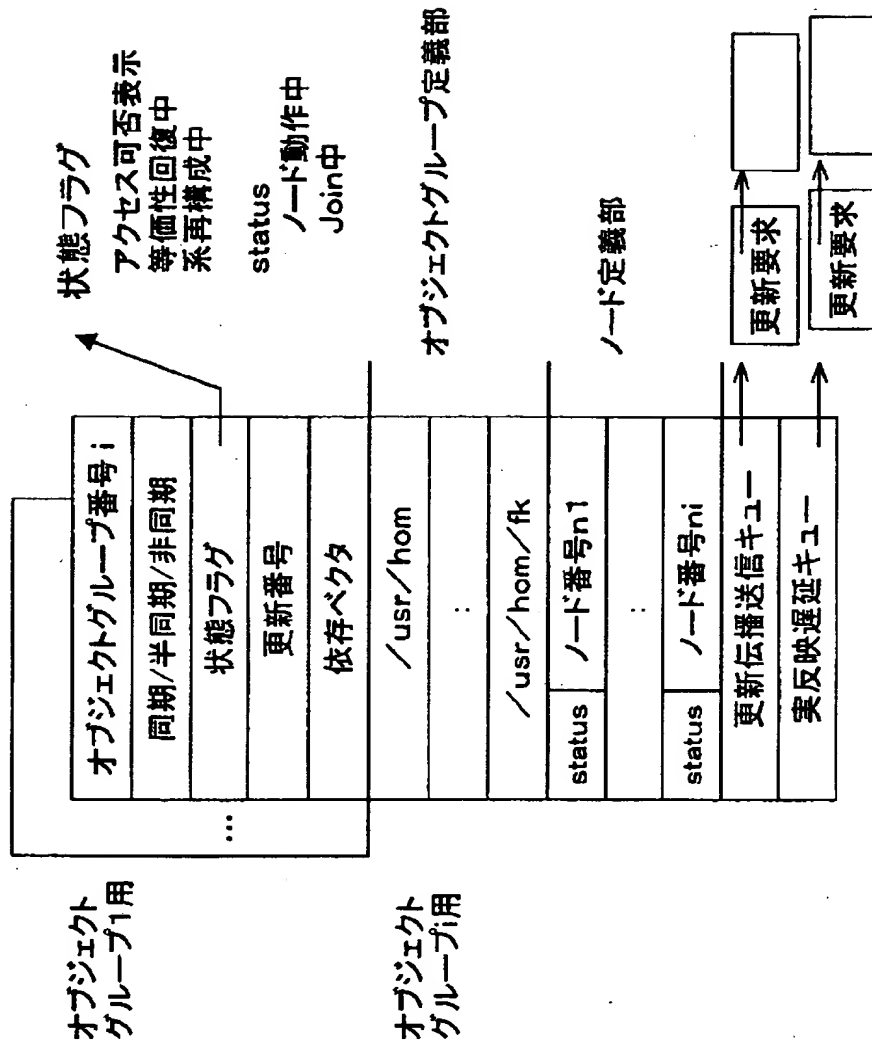
【図5】

系状態テーブルの構成例を示す図



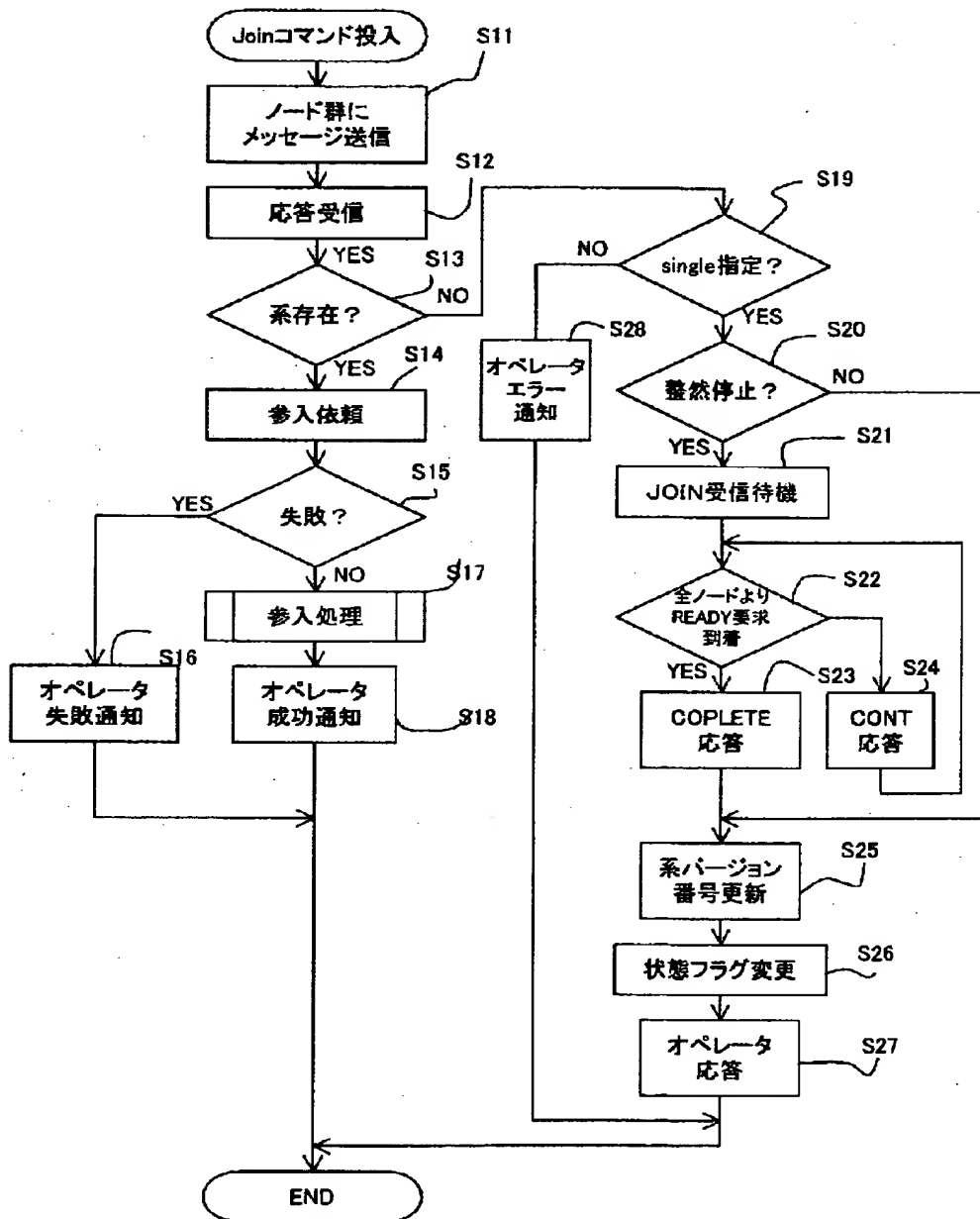
【図 6】

内部制御表の構成例を示す図



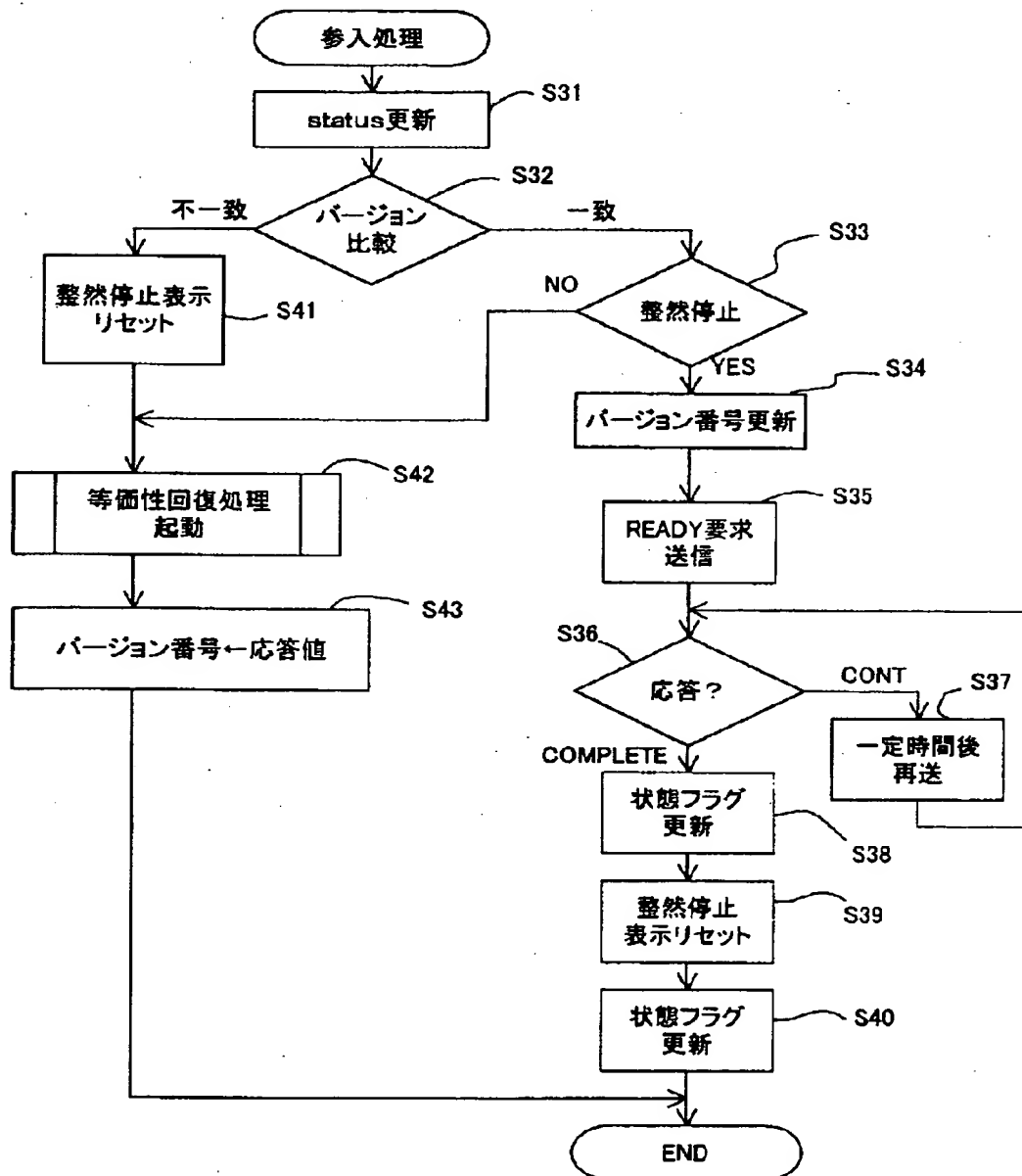
【図 7】

Joinコマンド投入時の系構成管理部による
動作処理を示すフローチャート



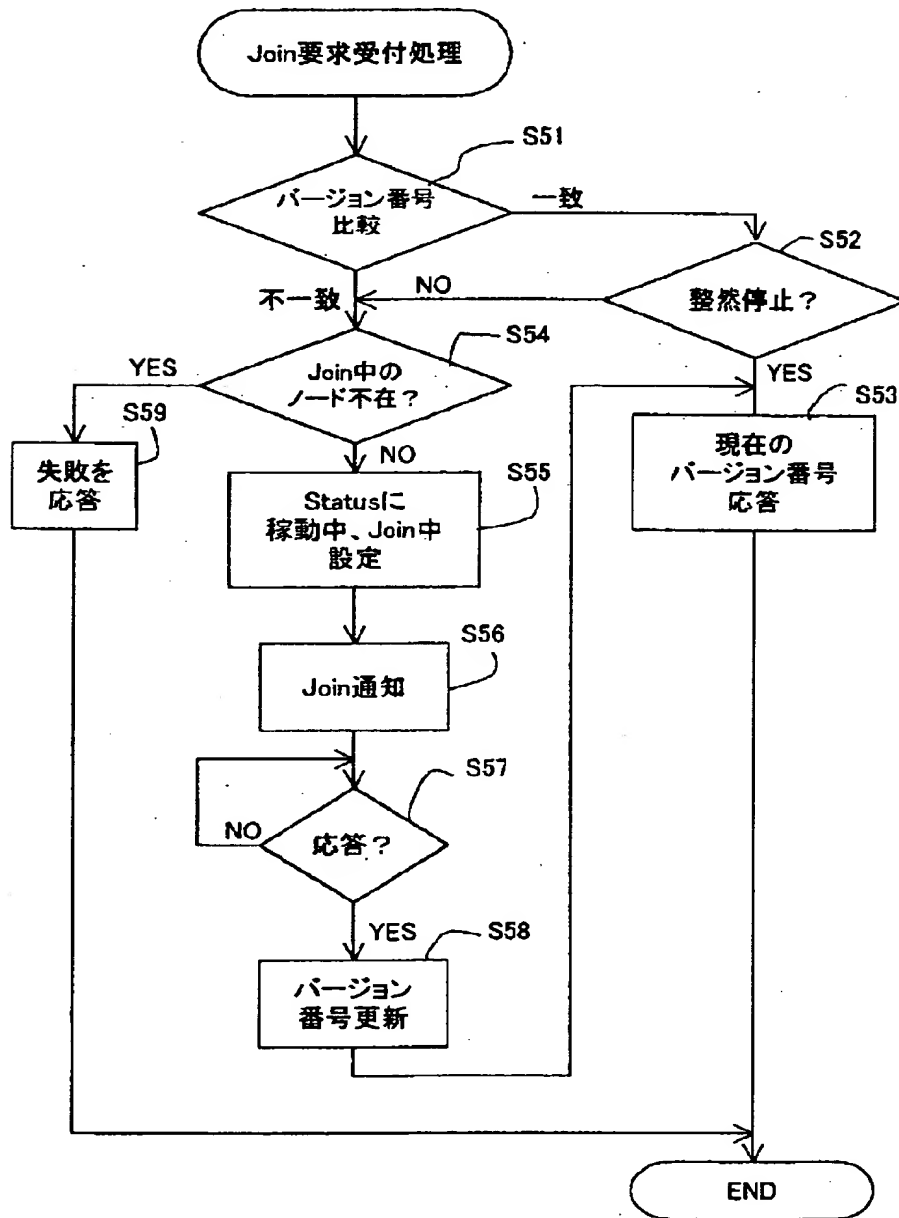
【図 8】

参入処理時の系構成管理部の動作処理
を示すフローチャート



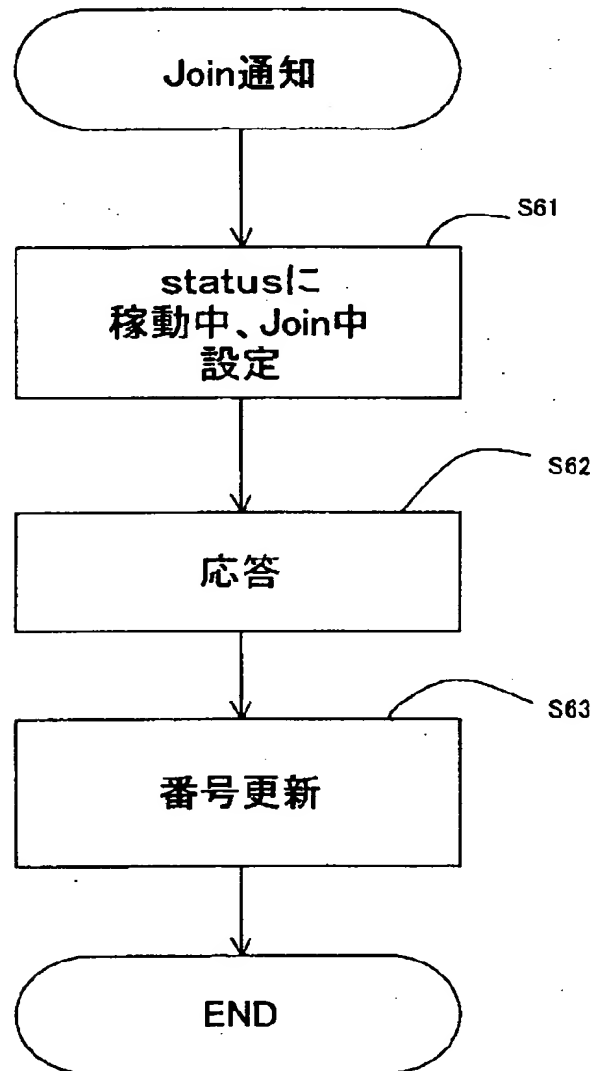
【図9】

JOIN要求受付処理時の系構成管理部の
動作処理を示すフローチャート



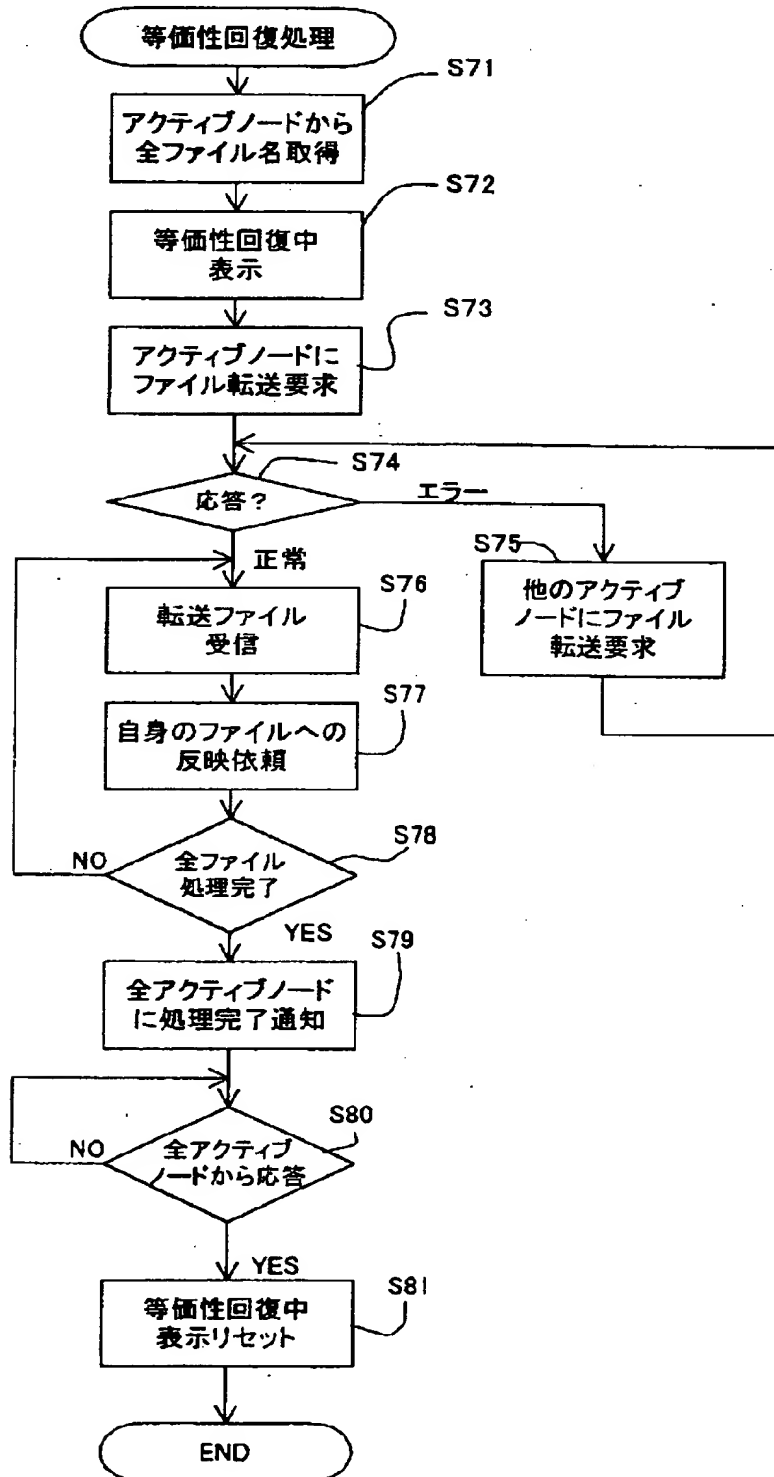
【図 1 0】

Join通知を受取ったノードの系構成管理部が
行う処理を示すフローチャート



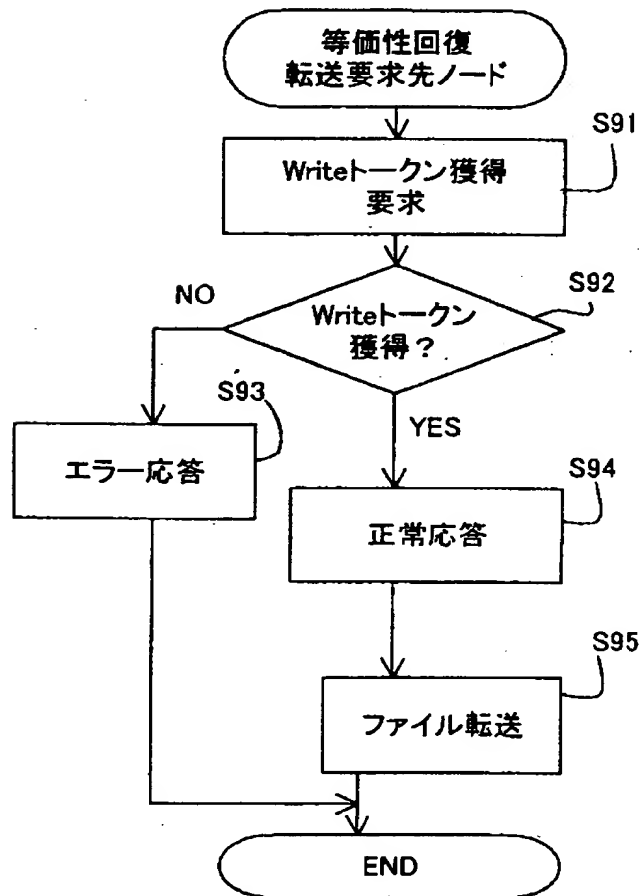
【図 1 1】

等価性回復処理の系構成管理部の 動作処理を示すフローチャート



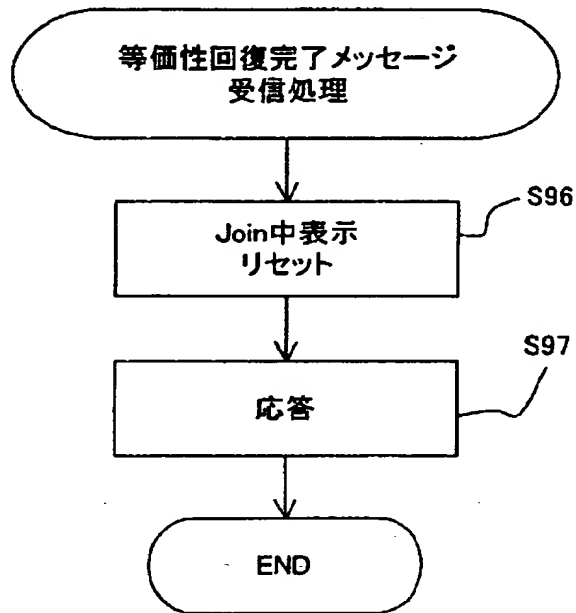
【図 1 2】

等価性回復転送要求を受信した
ノードの系構成管理部が行う
動作処理を示すフローチャート



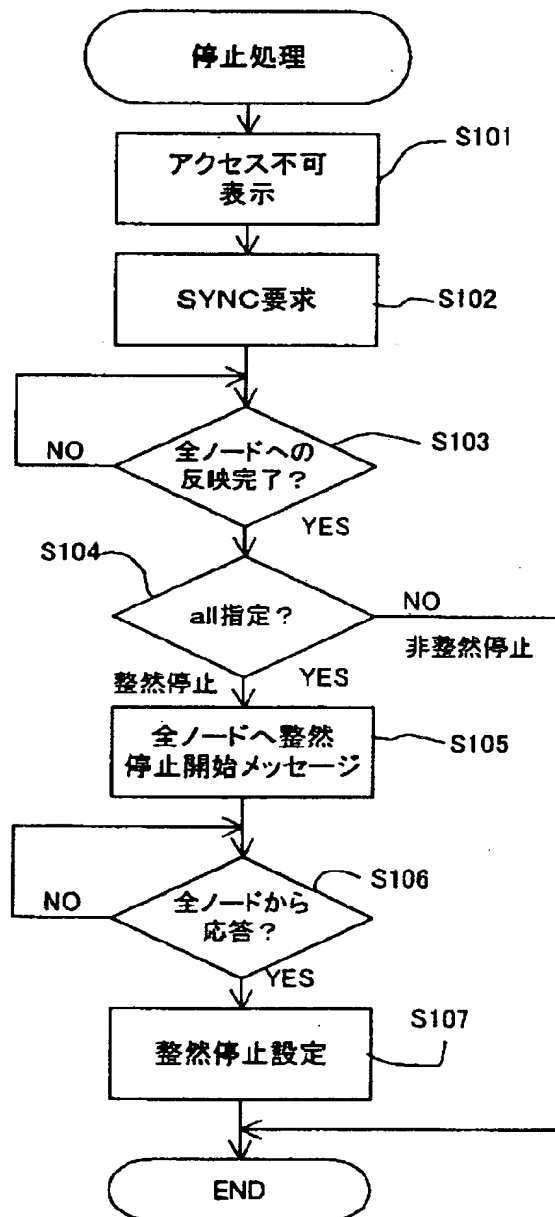
【図 1 3】

等価性回復完了メッセージを
受信したノードの系構成管理部
が行う動作処理を示すフローチャート



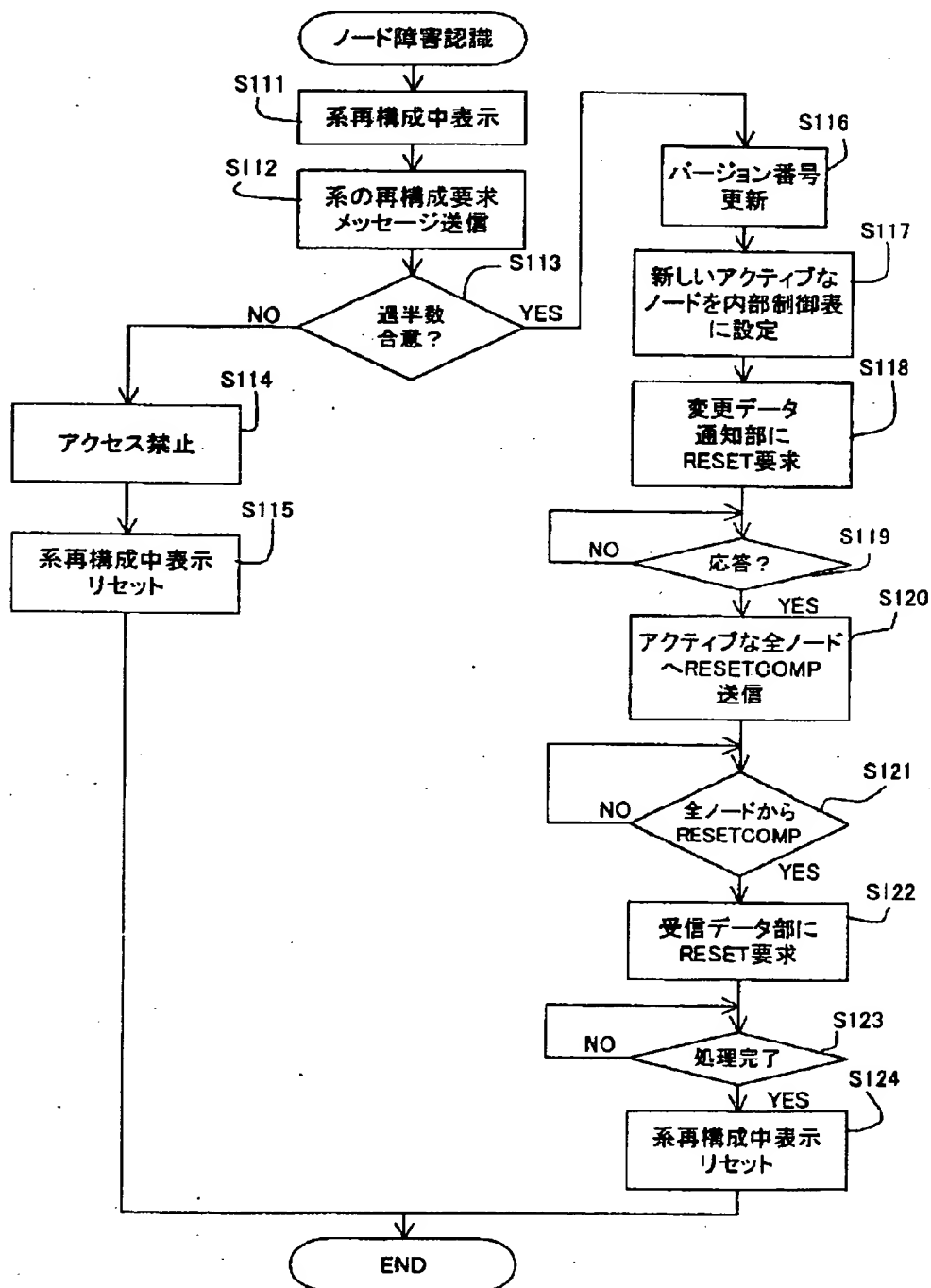
【図 1 4】

leaveコマンドを投入された時の
系構成管理部の動作処理を示すフローチャート



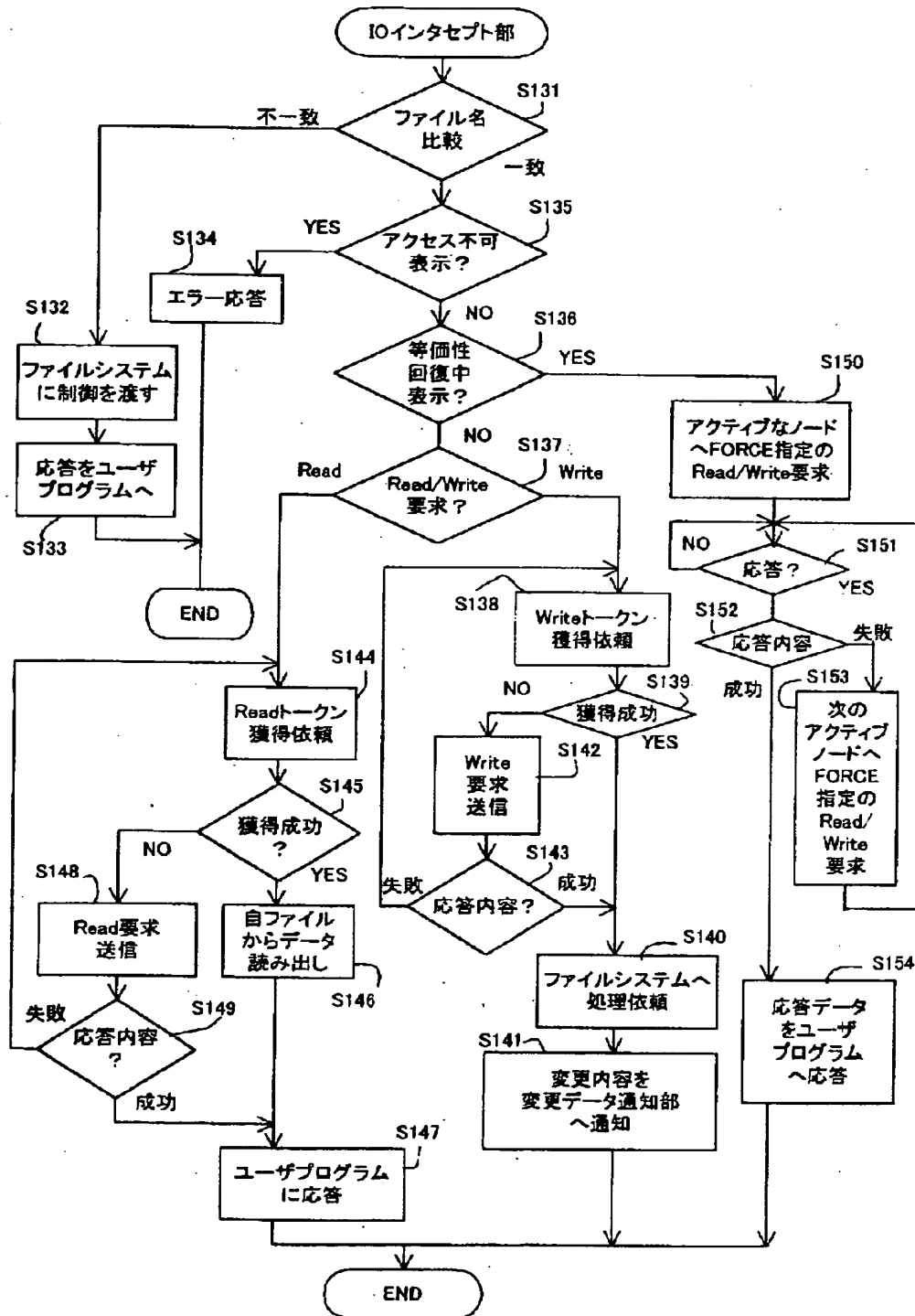
【図 15】

系内の他ノードの離脱を認識したノードの
系構成管理部の処理動作を示すフローチャート



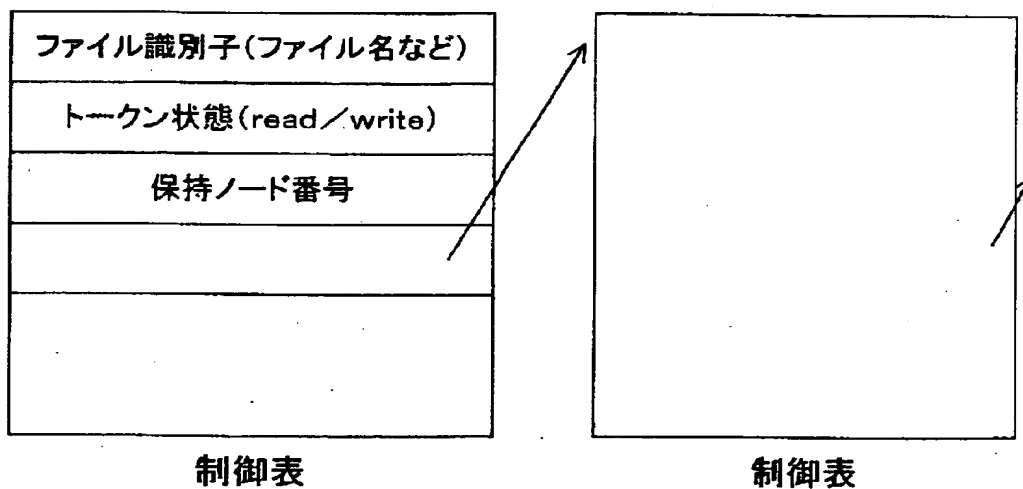
【図 16】

IO要求インタセプト部による処理動作を示すフローチャート



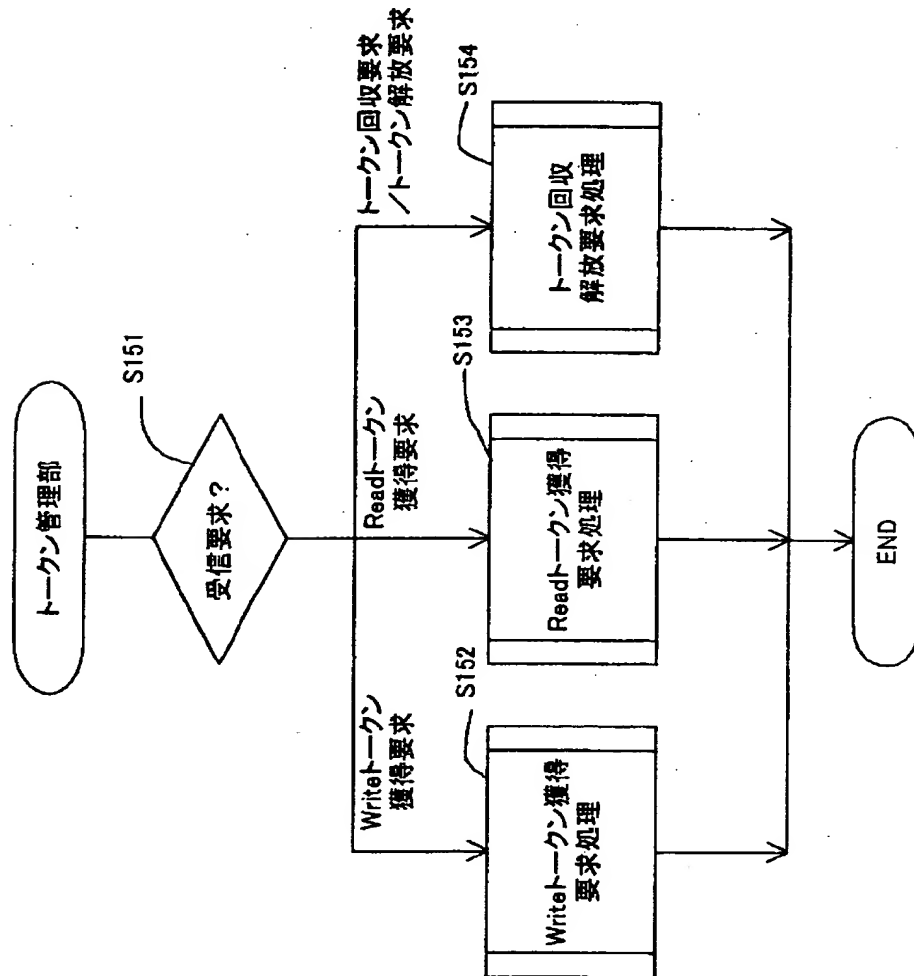
【図 1 7】

トークン制御表の構成例を示す図



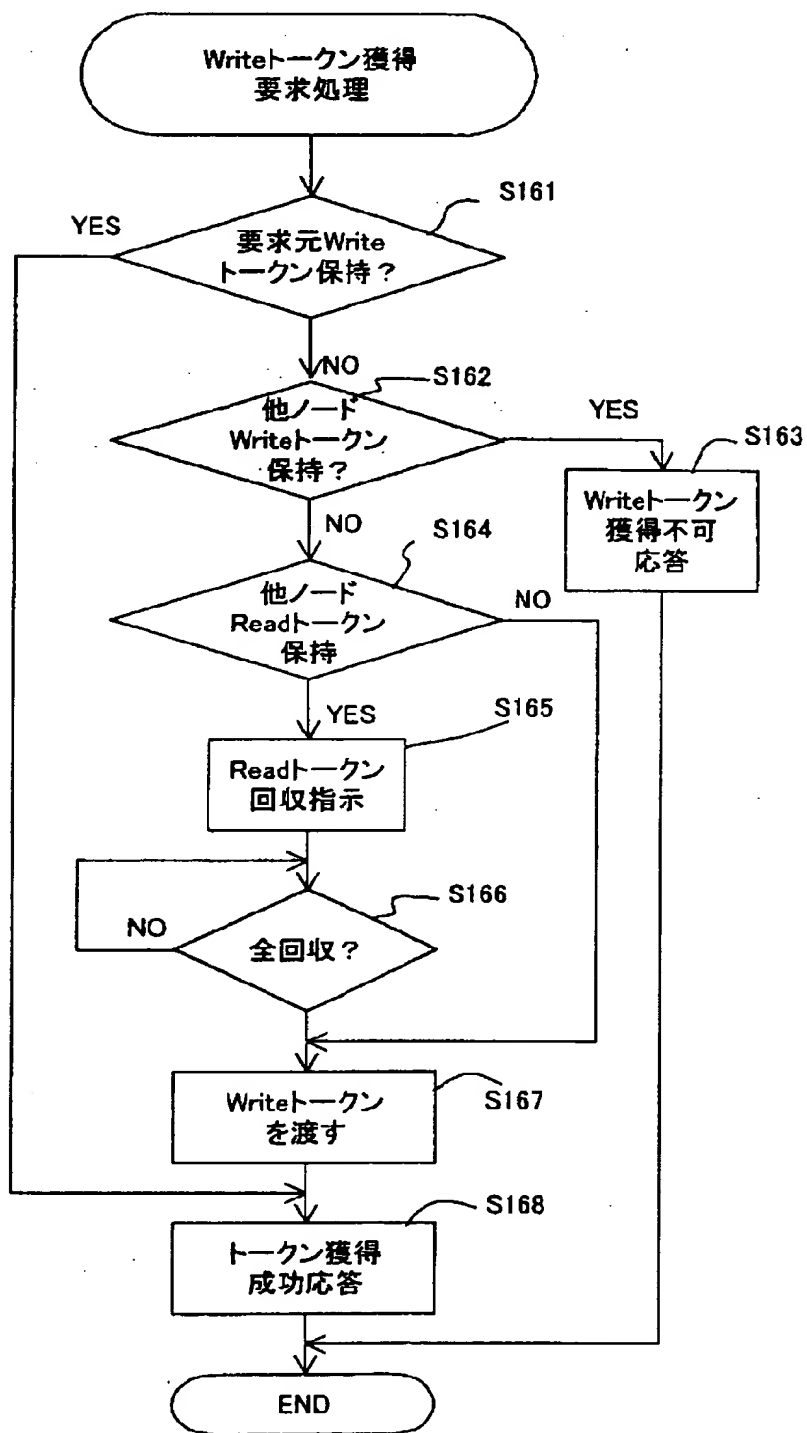
【図 1 8】

トークン管理部ノードのトークン管理部
の処理動作を示すフローチャートである



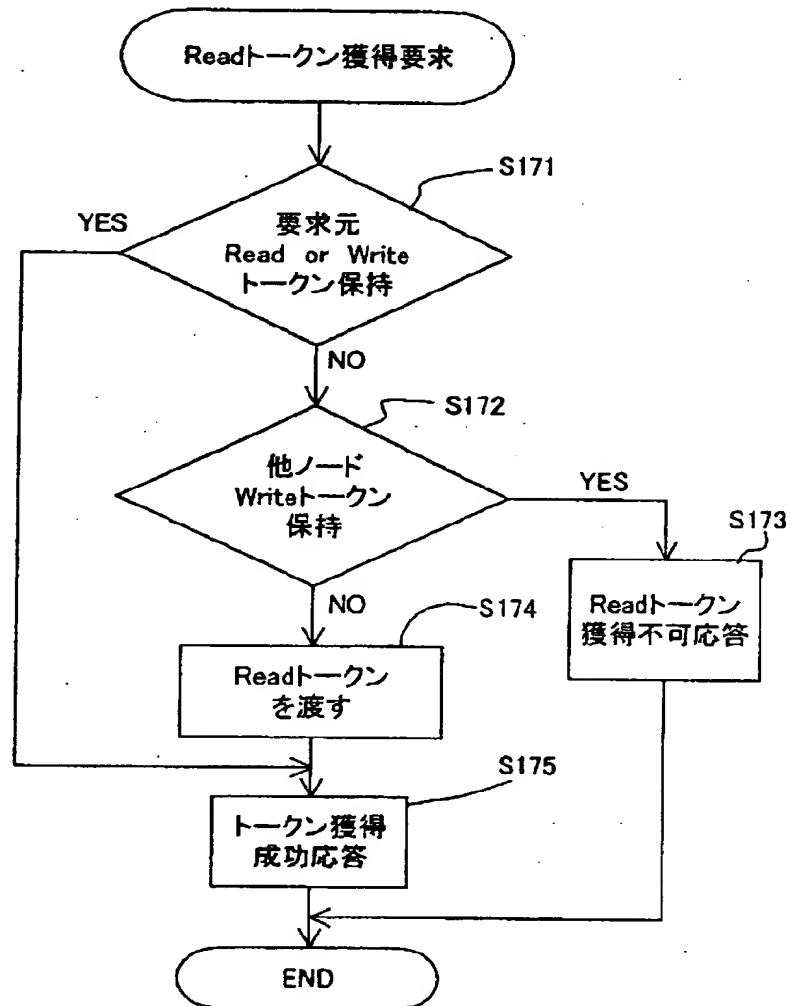
【図 1 9】

Writeトークン獲得要求処理時のトークン管理部 の処理動作を示すフローチャート



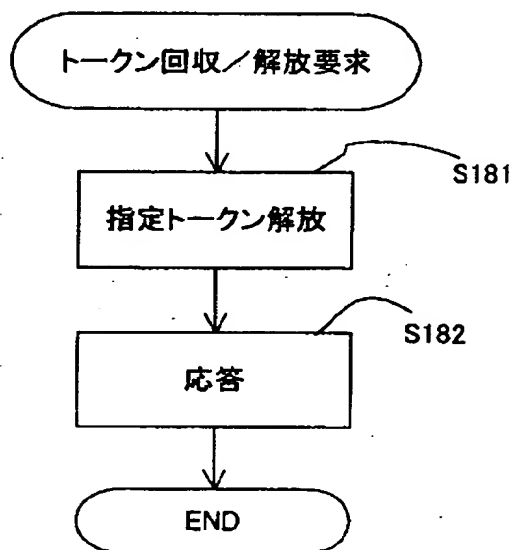
【図 2 0】

Readトークン獲得要求処理時のトークン管理部
の処理動作を示すフローチャート



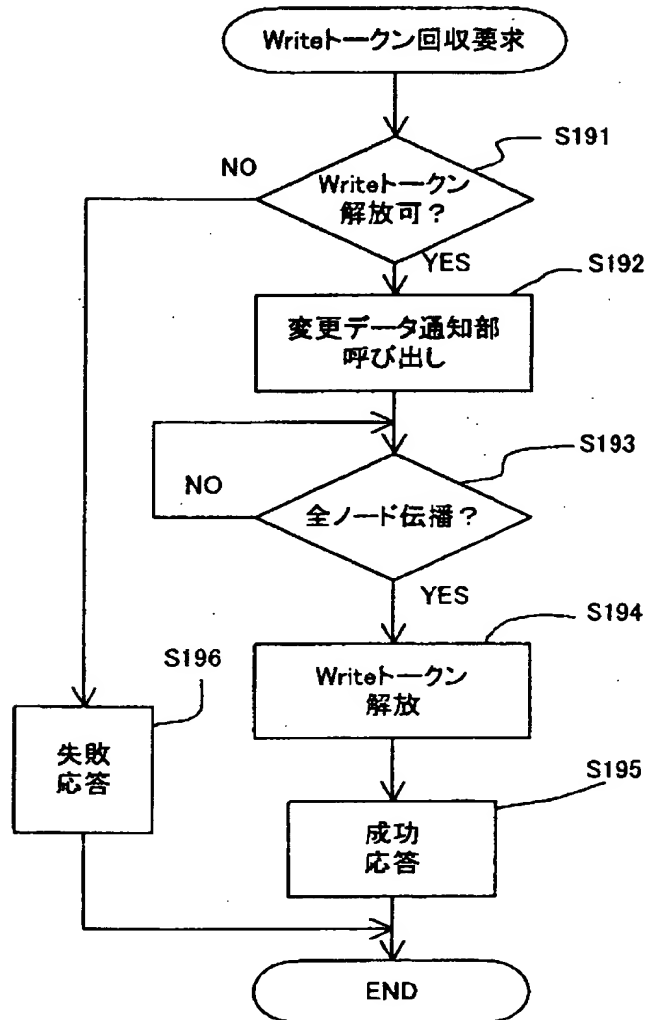
【図 2 1】

トークン解放／回収要求処理時
のトークン管理部の処理動作を
示すフローチャート



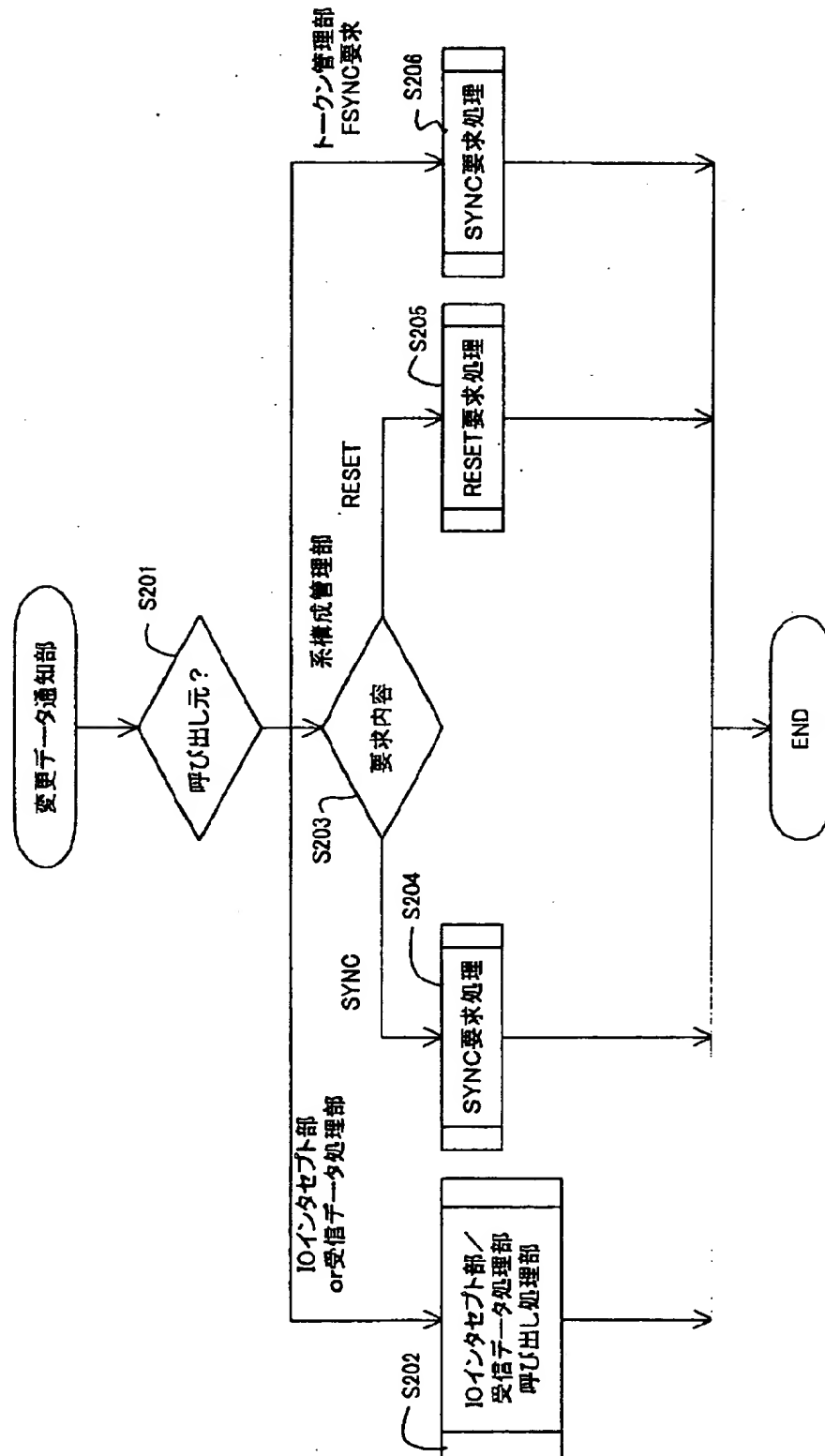
【図 2 2】

Writeトークン回収要求を受けたWriteトークン
保持ノードが行う動作処理を示すフローチャート



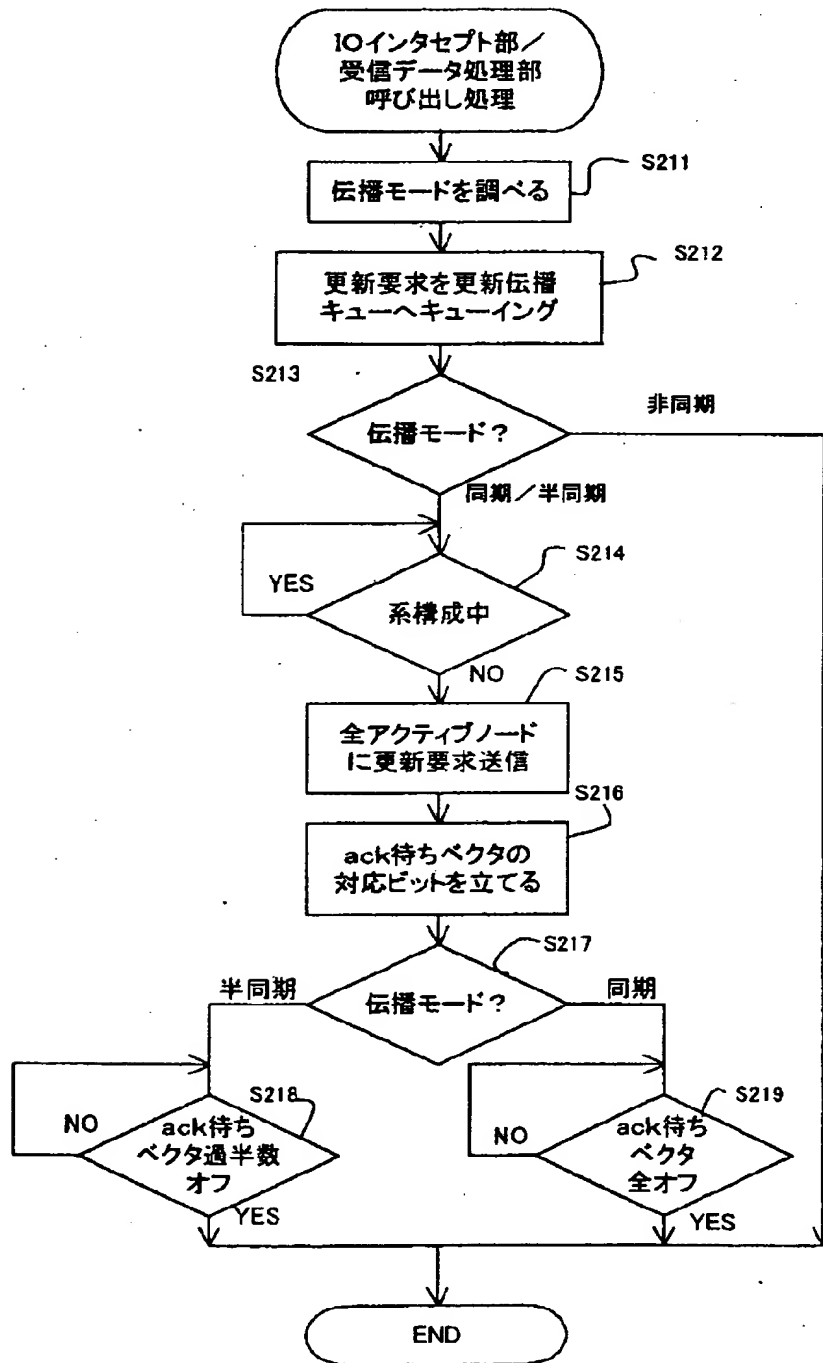
【図 2 3】

変更データ通知部による動作処理を示すフローチャート



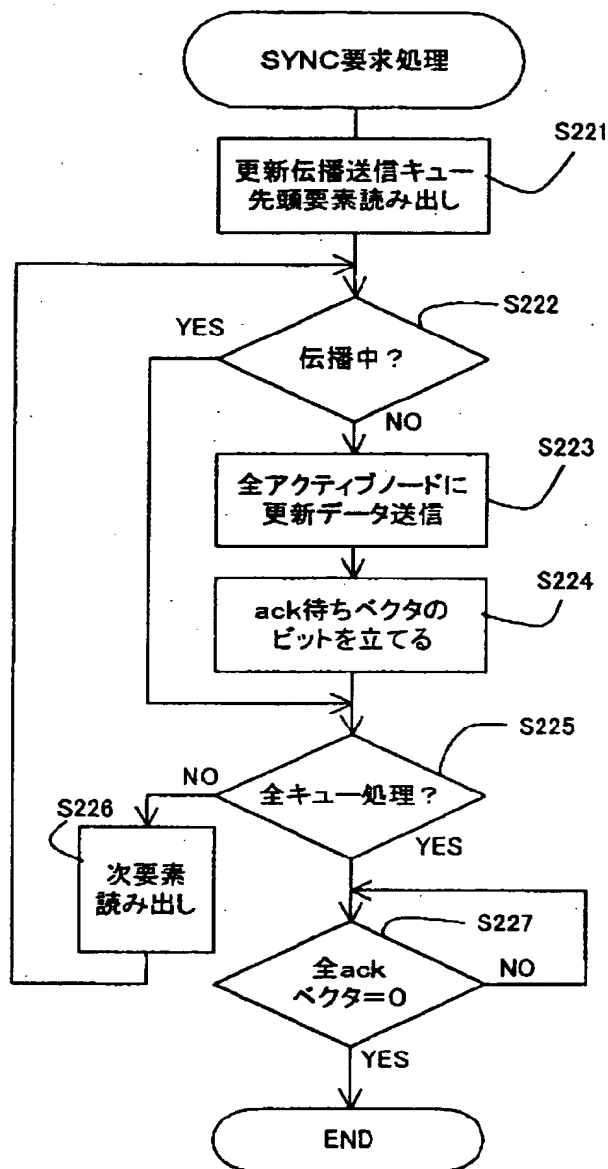
【図 2 4】

IO要求インタセプト部／受信データ処理部
呼び出し処理の変更データ通知部
の動作処理を示すフローチャート



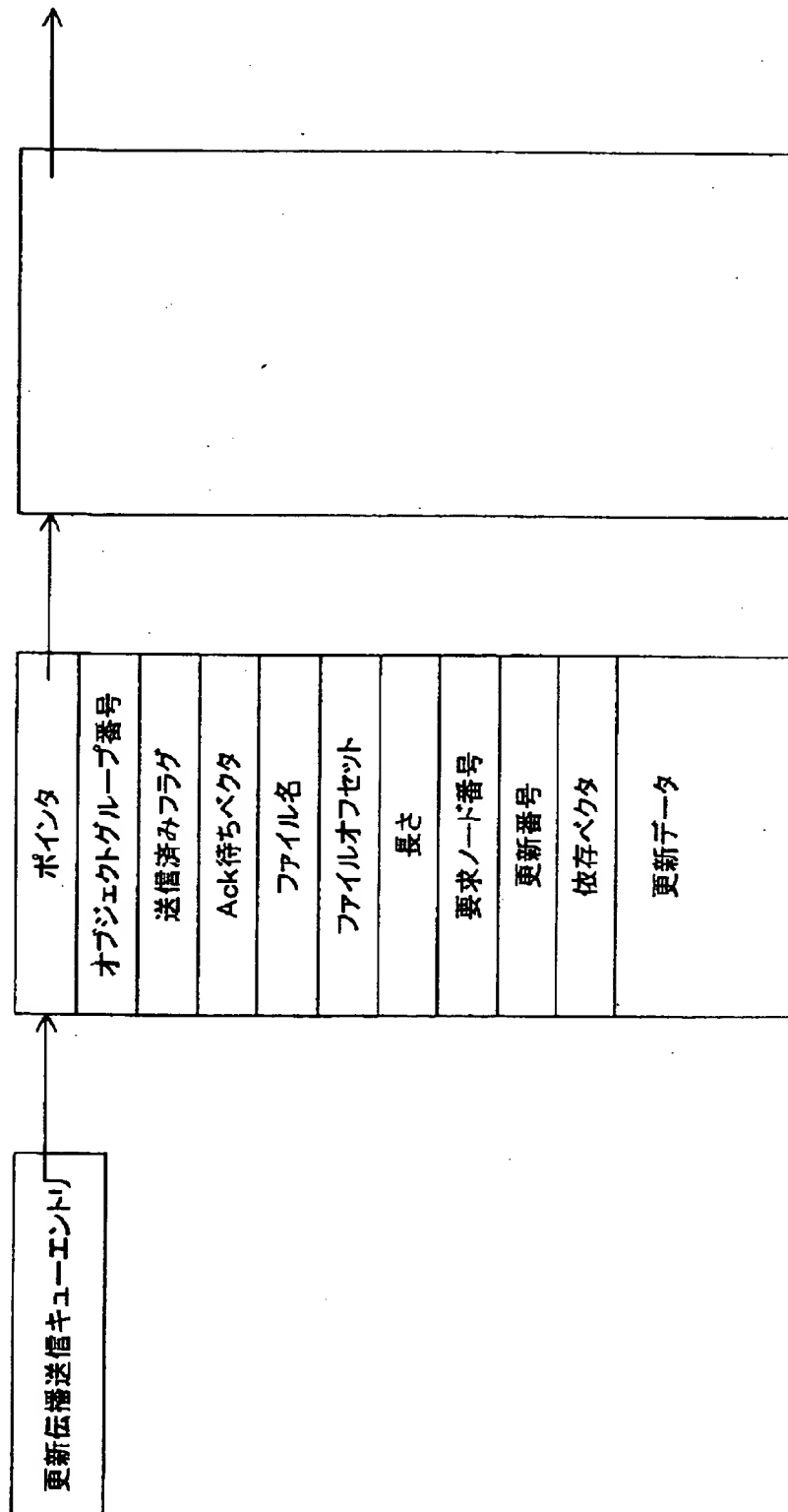
【図 2 5】

SYNC要求処理時の変更データ通知部の 動作処理を示すフローチャート



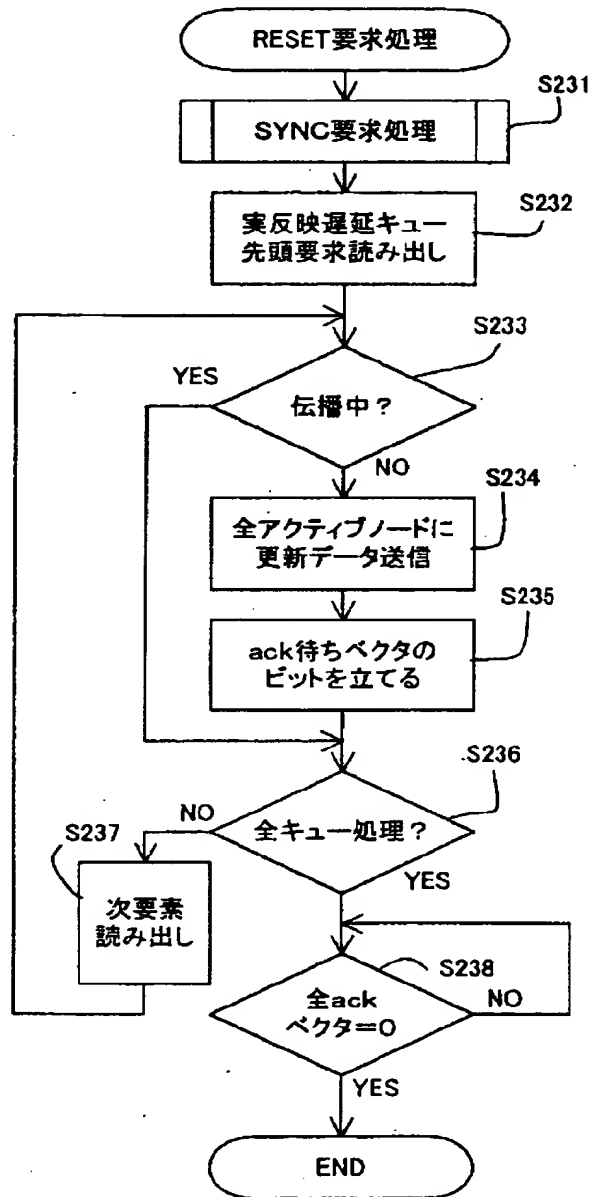
【図 2 6】

更新伝播送信キューの構成例を示す図



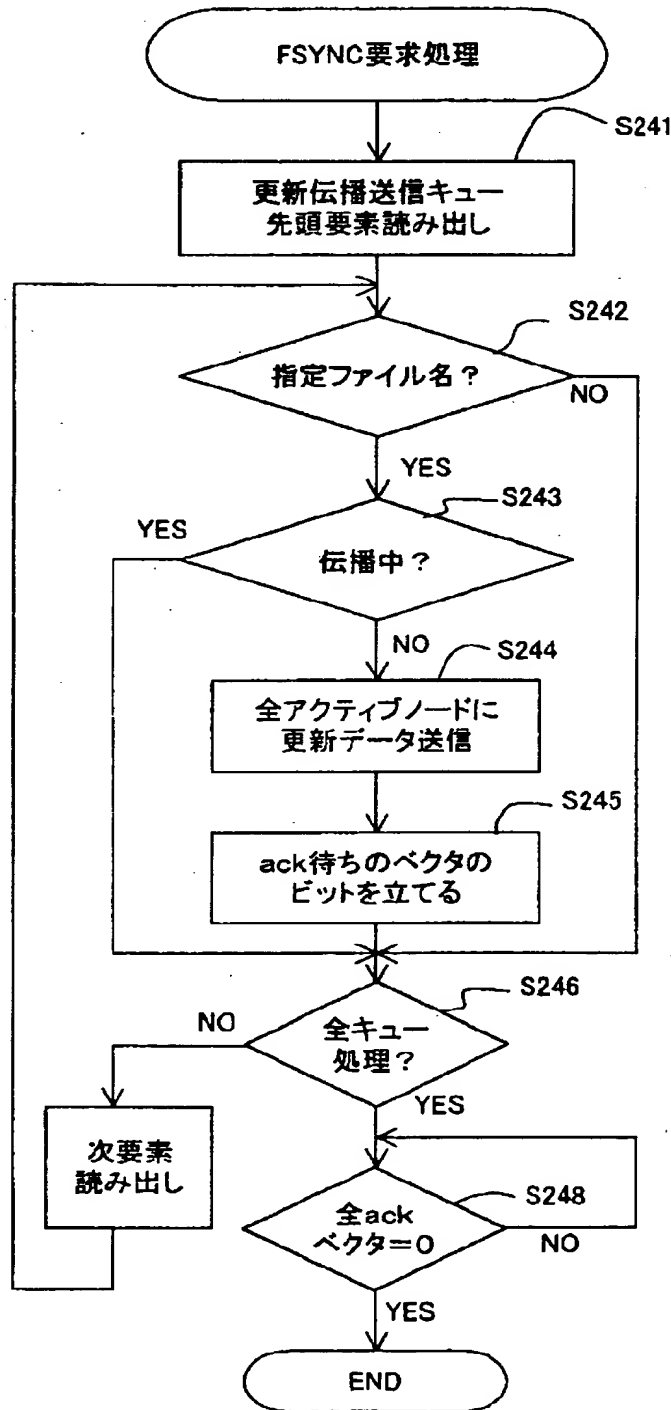
【図 2 7】

RESET要求処理時の変更データ通知部
の動作処理を示すフローチャート



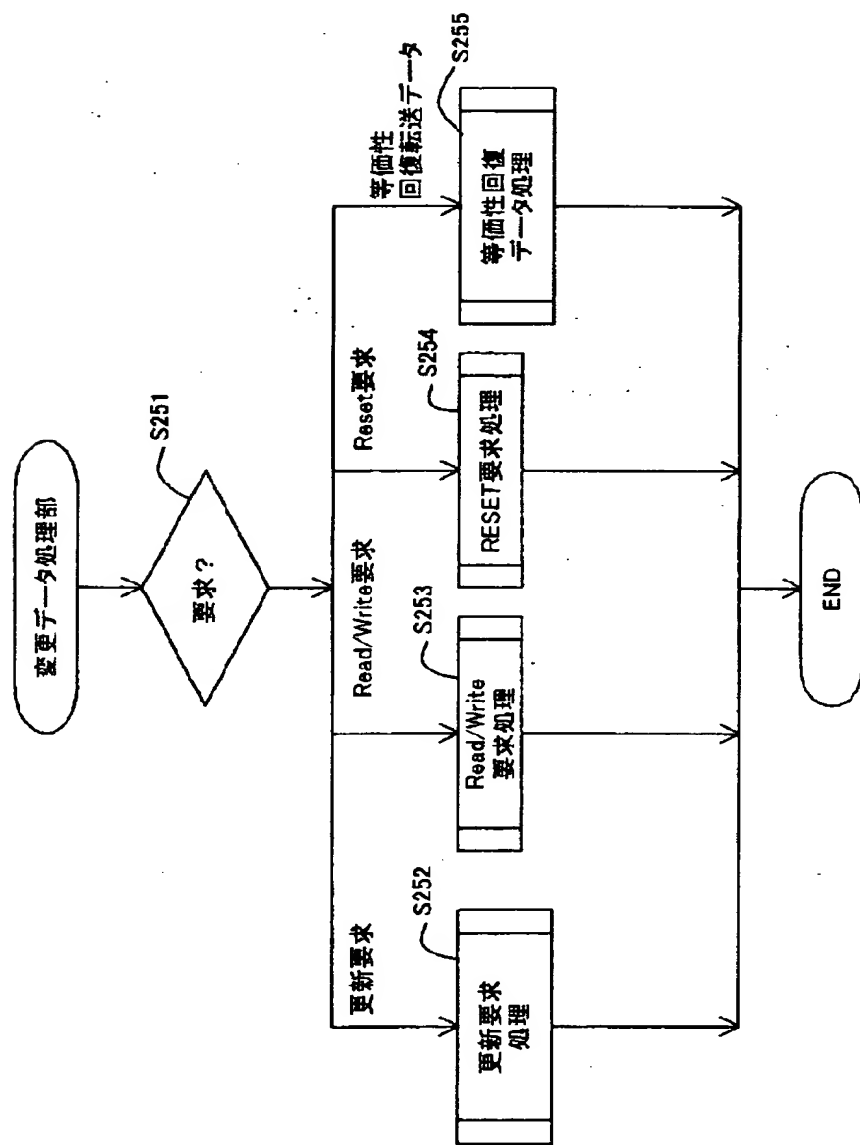
【図 28】

FSYNC要求処理時の変更データ通知部の
動作処理を示すフローチャート



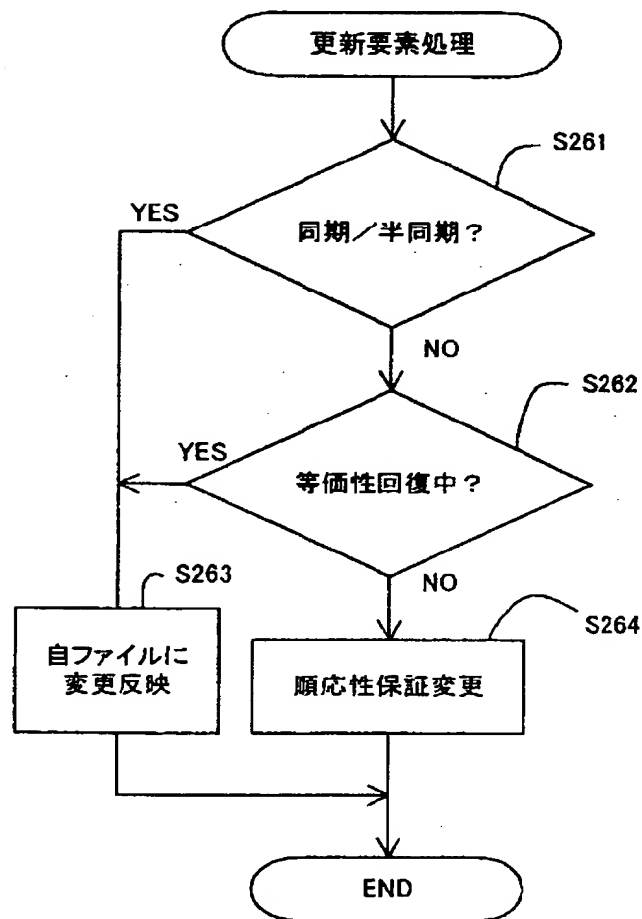
【図 2 9】

受信データ処理部の動作処理を示すフローチャート



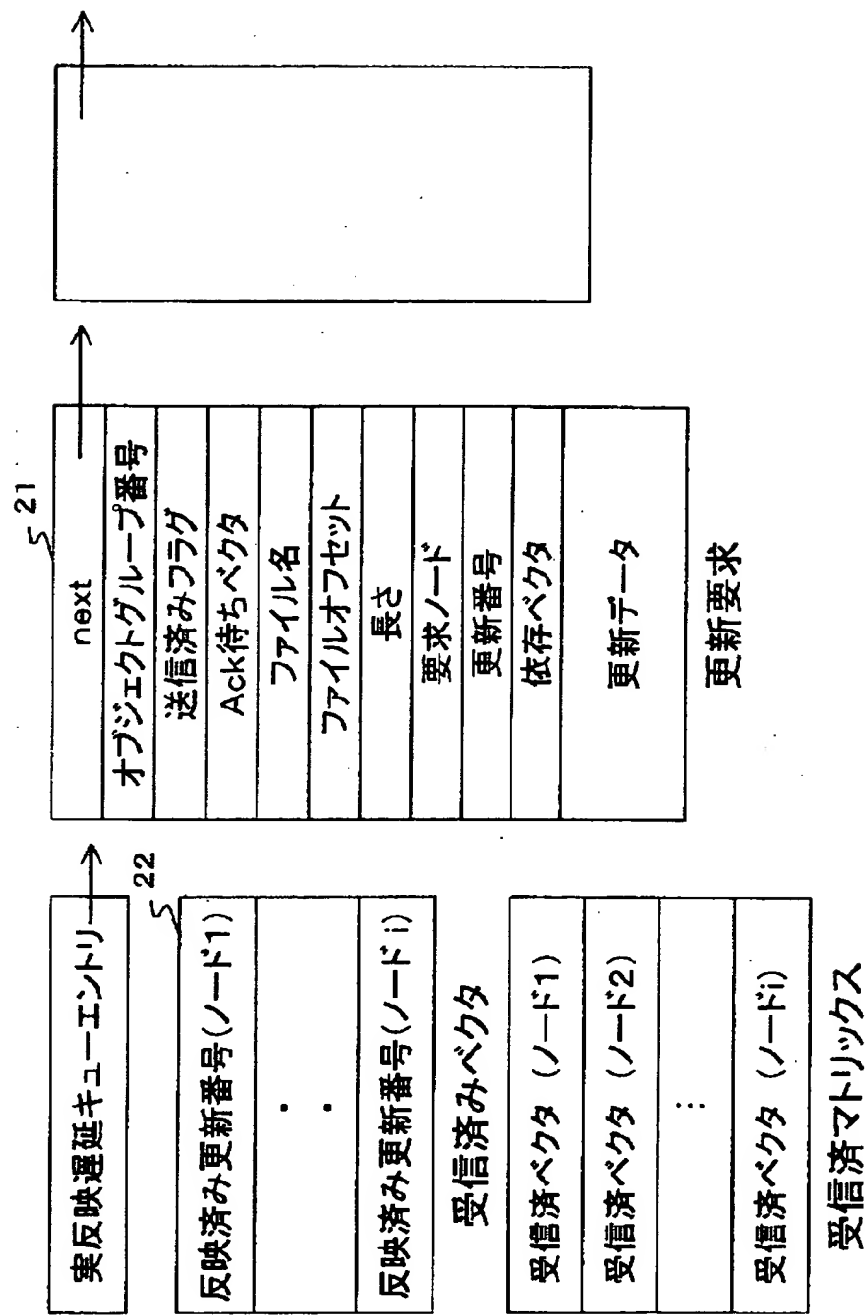
【図 3 0】

更新要求処理における受信データ処理部の
動作処理を示すフローチャート



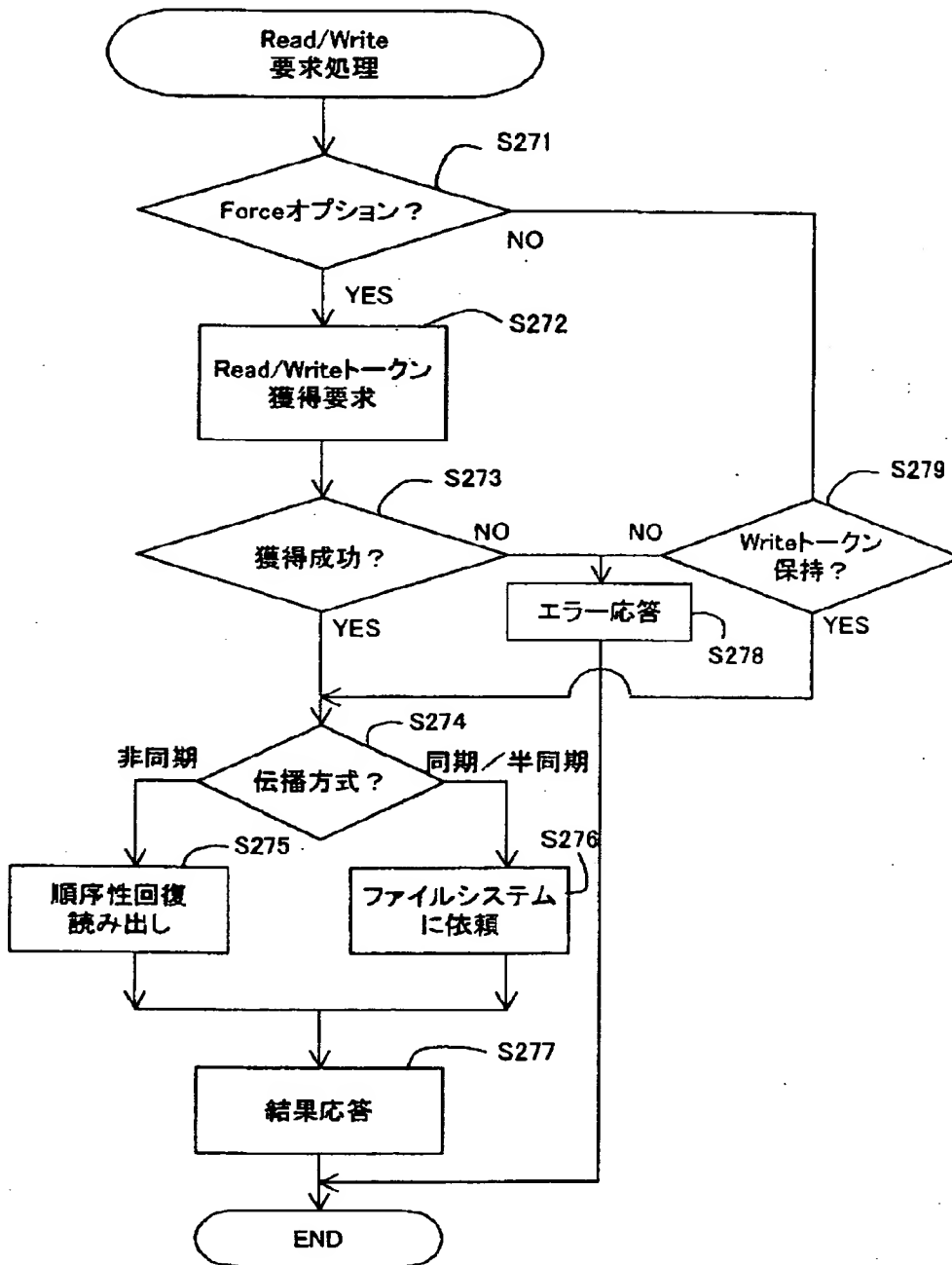
【図 3 1】

実反映遅延キューの構成例を示す図



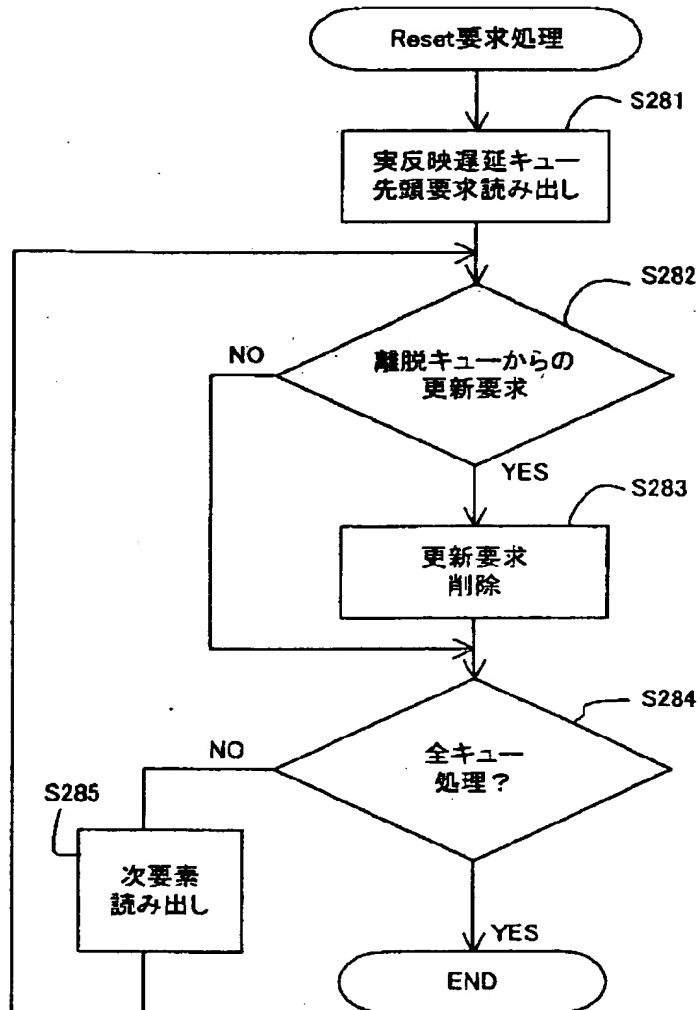
【図 3 2】

Read/Write要求処理における受信データ処理部の
の処理を示すフローチャート



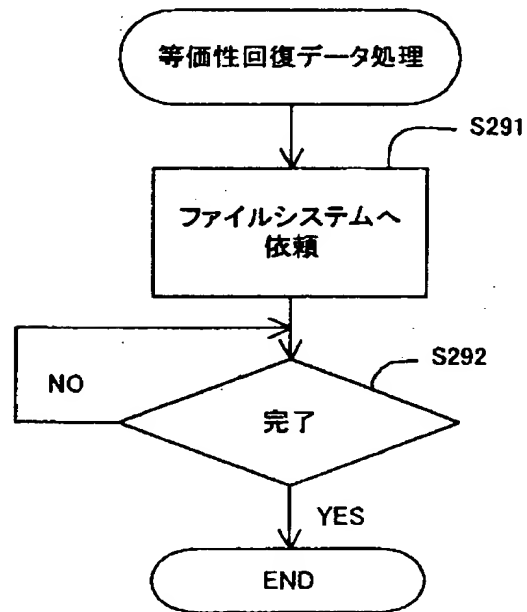
【図 3 3】

Reset要求処理における受信データ処理部の
動作処理を示すフローチャート



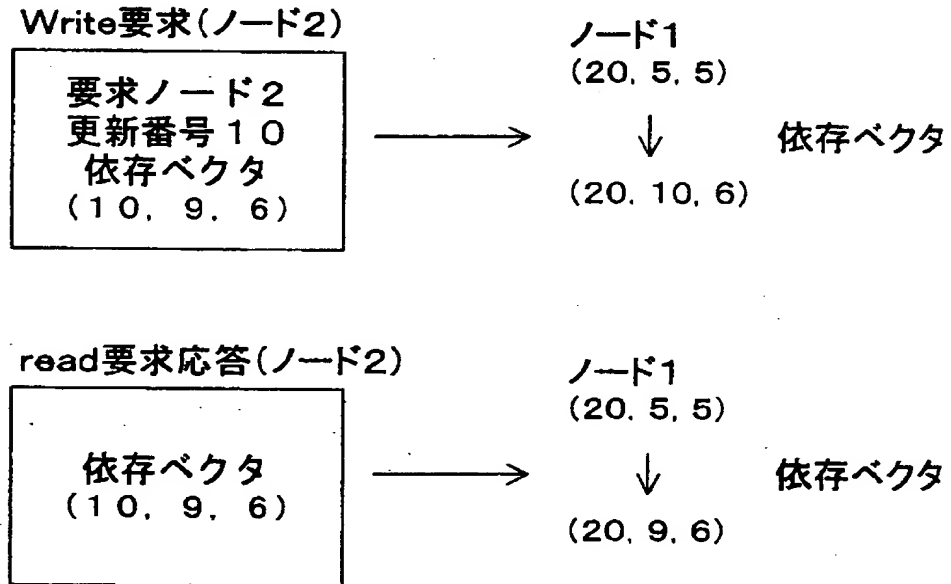
【図 3 4】

等価性回復データ処理における
受信データ処理部の動作処理
を示すフローチャート



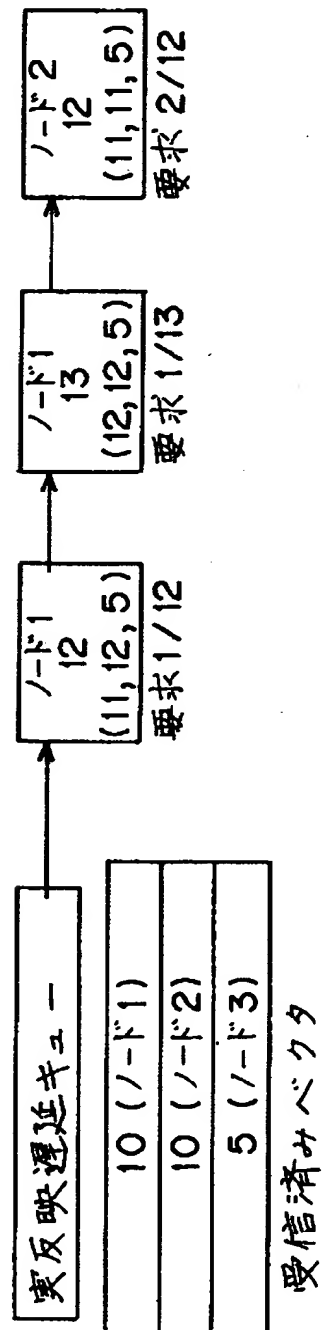
【図 3 5】

Write要求及びRead要求の応答に付加される
依存ベクタの例を示す図



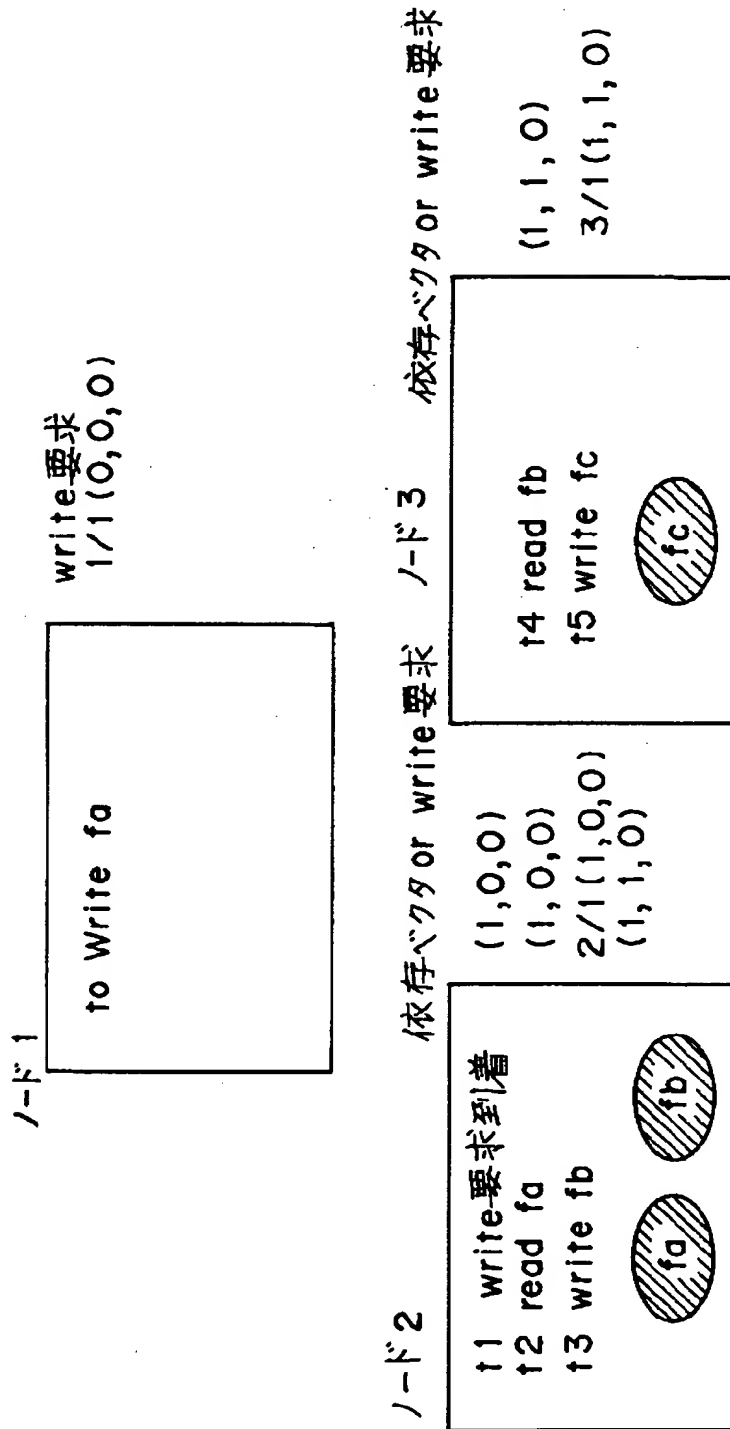
【図 36】

Write 要求及び Read 要求の応答に
付加される依存ベクタの例を示す図



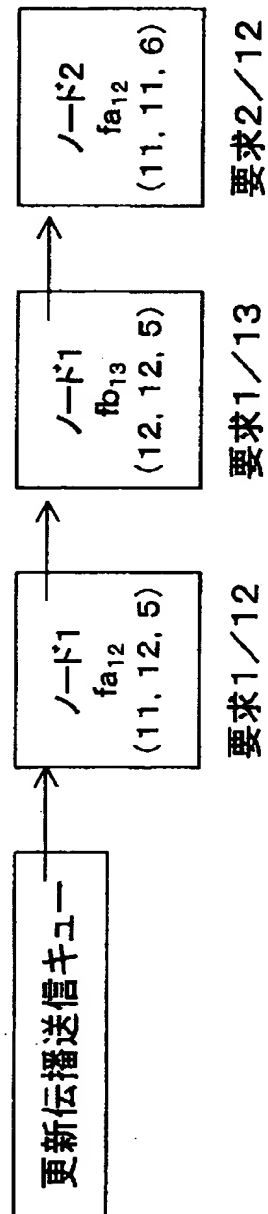
【図 3 7】

依存関係のある更新要求の順序性の保証を示す図



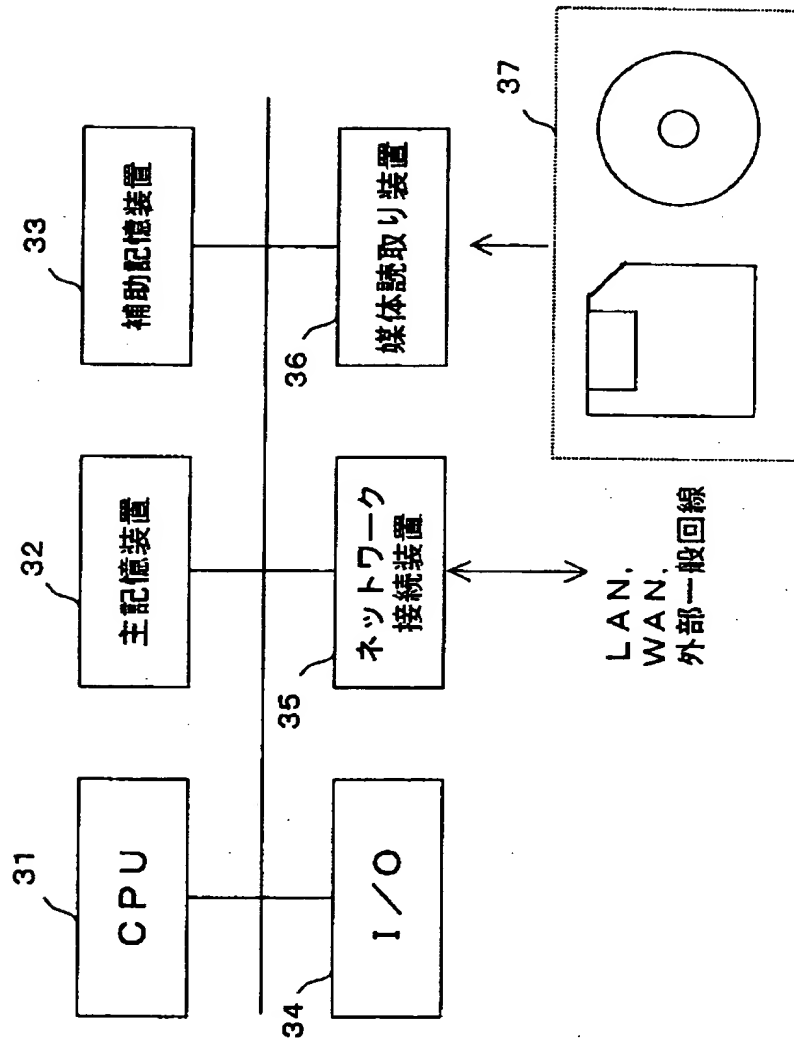
【図 3 8】

他ノードからのWrite要求を処理する時に、
更新伝播送信キューに同じファイルに対する
更新要求が存在していた場合の処理を説明する図



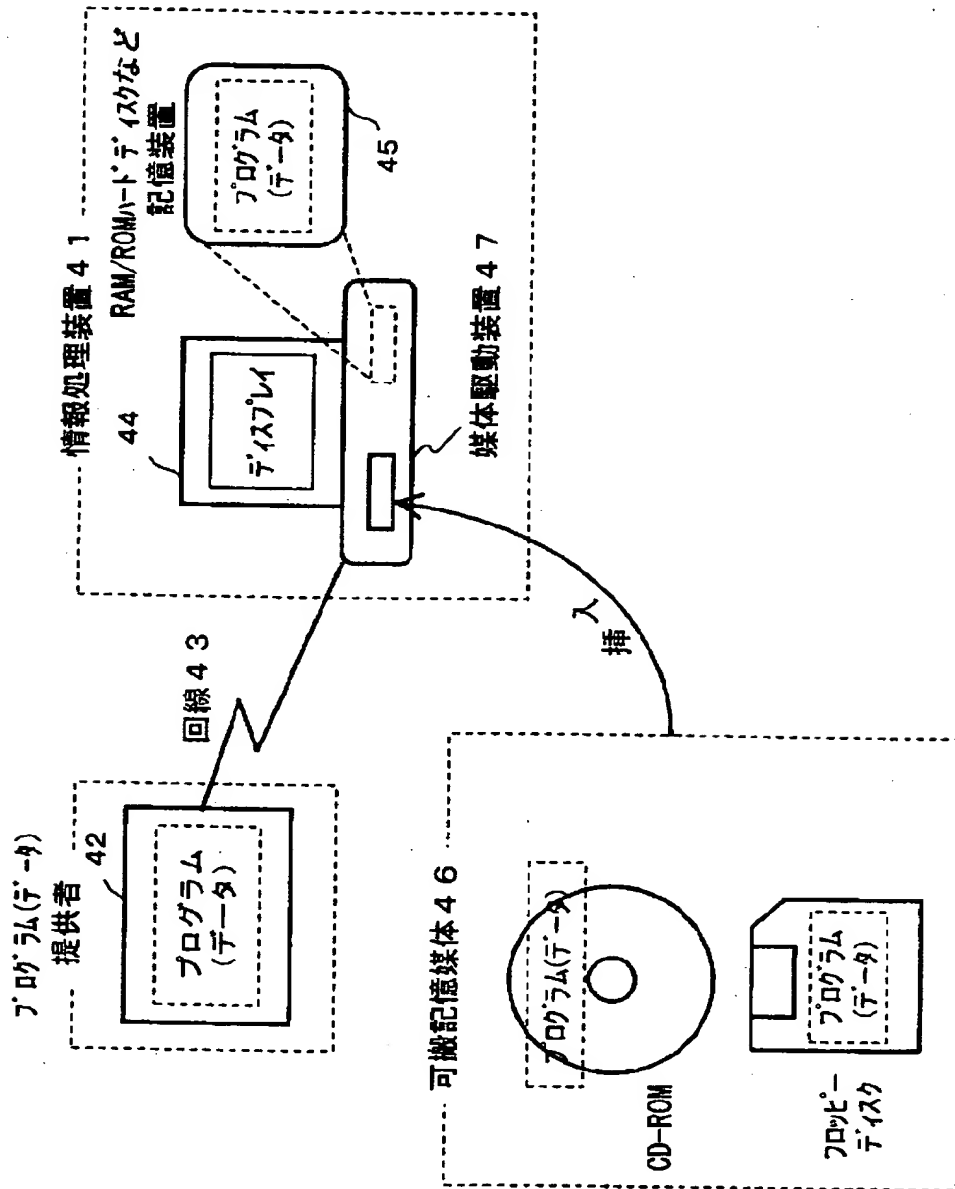
【図 39】

ノードとなる計算機のシステムの環境図



【図 40】

記 憶 媒 体 の 例 を 示 す 図



【書類名】 要約書

【要約】

【課題】 一貫性保証制御を付加し、ファイルレプリケーションの適用範囲を拡大したファイルレプリケーションシステムを提供することを課題とする。

【解決手段】 自ノード内で生じた共用ファイル 6 に対するアクセス要求に対し、I/O 要求インタセプト手段 2 はトークン管理手段 3 に共用ファイル 6 に対するアクセス許可を求める。トークン管理手段 6 は、トークン管理手段 1 3 は、I/O 要求インタセプト手段 2 からのアクセス許可要求に対し、既に他のノードが共用ファイルに対する更新許可を保持する時、更新許可を保持するノードを I/O 要求インタセプト手段 2 に通知する。I/O 要求インタセプト手段は、トークン管理手段からアクセス許可が得られない時、通知された更新許可を保持するノードに共用ファイルへのアクセス処理を依頼する。

【選択図】 図 1

出 願 人 履 歴 情 報

識別番号 [000005223]

1. 変更年月日	1996年 3月26日
[変更理由]	住所変更
住 所	神奈川県川崎市中原区上小田中4丁目1番1号
氏 名	富士通株式会社